

SANDIA REPORT

SAND2022-8810

Unclassified Unlimited Release

Printed June 2022



Sandia
National
Laboratories

The Portals 4.3 Network Programming Interface

Brian W. Barrett, Ron Brightwell, Ryan E. Grant, Whit Schonbein, Scott Hemmert,
Kevin Pedretti, Keith Underwood, Rolf Riesen, Torsten Hoefler, Mathieu Barbe, Luiz H. Suraty Filho,
Alexandre Ratchov, and Arthur B. Maccabe

Prepared by
Sandia National Laboratories
Albuquerque, New Mexico 87185
Livermore, California 94550

Issued by Sandia National Laboratories, operated for the United States Department of Energy by National Technology & Engineering Solutions of Sandia, LLC.

NOTICE: This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government, nor any agency thereof, nor any of their employees, nor any of their contractors, subcontractors, or their employees, make any warranty, express or implied, or assume any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represent that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government, any agency thereof, or any of their contractors or subcontractors. The views and opinions expressed herein do not necessarily state or reflect those of the United States Government, any agency thereof, or any of their contractors.

Printed in the United States of America. This report has been reproduced directly from the best available copy.

Available to DOE and DOE contractors from

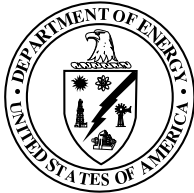
U.S. Department of Energy
Office of Scientific and Technical Information
P.O. Box 62
Oak Ridge, TN 37831

Telephone: (865) 576-8401
Facsimile: (865) 576-5728
E-Mail: reports@osti.gov
Online ordering: <http://www.osti.gov/scitech>

Available to the public from

U.S. Department of Commerce
National Technical Information Service
5301 Shawnee Road
Alexandria, VA 22312

Telephone: (800) 553-6847
Facsimile: (703) 605-6900
E-Mail: orders@ntis.gov
Online order: <https://classic.ntis.gov/help/order-methods>



The Portals 4.3 Network Programming Interface

Ron Brightwell
Whit Schonbein
Kevin Pedretti
Scalable System Software Department

Scott Hemmert
Scalable Computer Architecture Department

Sandia National Laboratories
P.O. Box 5800
Albuquerque, NM 87185-1319
{rbbrigh, wwschon,
ktpedre, kshemme}@sandia.gov

Arthur B. Maccabe
Computer Science and Mathematics
Oak Ridge National Laboratory
Oak Ridge, TN 37831
maccabeab@ornl.gov

Ryan E. Grant
Queen's University
Kingston, Ontario, Canada
ryan.grant@queensu.edu

Brian W. Barrett
Amazon.com, Inc.
brian@bbarrett.org

Keith Underwood
HPE
keith.underwood@hpe.com

Rolf Riesen
Intel Corporation
rolf.riesen@intel.com

Torsten Hoefler
Computer Science Department
ETH Zurich
htor@inf.ethz.ch

Mathieu Barbe
Luiz Henrique Suraty Filho
Alexandre Ratchov
Atos
mathieu.barbe@atos.net
luiz@suraty.com
alex@caoua.org

ABSTRACT

This report presents a specification for the Portals 4 network programming interface. Portals 4 is intended to allow scalable, high-performance network communication between nodes of a parallel computing system. Portals 4 is well suited to massively parallel processing and embedded systems. Portals 4 represents an adaptation of the data movement layer developed for massively parallel processing platforms, such as the 4500-node Intel TeraFLOPS machine. Sandia's Cplant cluster project motivated the development of Version 3.0, which was later extended to Version 3.3 as part of the Cray Red Storm machine and XT line. Version 4 is targeted to the next generation of machines employing advanced network interface architectures that support enhanced offload capabilities.

Acknowledgments

Over the years, many people have helped shape, design, and develop Portals. We wish to thank: Eric Barton, Peter Braam, Jerrie Coffman, Lee Ann Fisk, David Greenberg, Eric Hoffman, Trammel Hudson, Gabi Istrail, Jeanette Johnston, Chu Jong, Clint Kaul, Roy Larsen, Mike Levenhagen, Kevin McCurley, Jim Otto, Bob Pearson, David Robboy, Mark Sears, Lance Shuler, Jim Schutt, Mack Stallcup, Todd Underwood, David van Dresser, Dena Vigil, Lee Ward, Kyle Wheeler, Stephen Wheat, and Frank Zago.

People who were influential in managing the project were: Bill Camp, Ed Barsis, Art Hale, and Neil Pundit

While we have tried to be comprehensive in our listing of the people involved, it is very likely that we have missed at least one important contributor. The omission is a reflection of our poor memories and not a reflection of the importance of their contributions. We apologize to the unnamed contributors.

Contents

List of Figures	11
List of Tables	12
List of Implementation Notes	13
Preface	14
Nomenclature	15
1. Introduction	17
1.1. Overview	17
1.2. Purpose	17
1.3. Background	18
1.4. Scalability	19
1.5. Communication Model	19
1.6. Zero Copy, OS Bypass, and Application Bypass	19
1.7. Faults	20
2. An Overview of the Portals API	21
2.1. Data Movement	21
2.2. Unreliable Datagrams	24
2.3. Usage	24
2.4. Completion Events	26
2.5. Portals Addressing	26
2.5.1. Lists and List Entries	28
2.5.2. Match Lists and Match List Entries	29
2.6. Modifying Data Buffers	30
2.7. Ordering	30
2.7.1. Short Message Ordering Semantics	32
2.7.2. Long Message Ordering Semantics	32
2.7.3. Relative Ordering of Operations in Overlapping Portals	32
2.7.4. Ordering of Unexpected Messages	33
2.7.5. Relaxing Message Ordering	33
2.8. Flow Control	33

2.9. Multi-Threaded Applications	34
3. The Portals API	35
3.1. Naming Conventions and Typeface Usage	35
3.2. Constants	35
3.2.1. Version Information	36
3.3. Base Types	36
3.3.1. Sizes	36
3.3.2. Handles	36
3.3.3. Indexes	37
3.3.4. Match Bits	37
3.3.5. Network Interfaces	37
3.3.6. Identifiers	37
3.3.7. Status Registers	37
3.4. Function Arguments and Return Codes	38
3.5. Initialization and Cleanup	38
3.5.1. PtlInit	39
3.5.2. PtlFini	39
3.5.3. PtlAbort	39
3.6. Network Interfaces	40
3.6.1. The Network Interface Limits Type	40
3.6.2. PtlNIInit	42
3.6.3. PtlNIFini	44
3.6.4. PtlNIStatus	44
3.6.5. PtlNIHandle	45
3.6.6. PtlSetMap	45
3.6.7. PtlGetMap	46
3.7. Portal Table Entries	47
3.7.1. PtlPTAlloc	47
3.7.2. PtlPTFree	49
3.7.3. PtlPTDisable	49
3.7.4. PtlPTEnable	50
3.8. Usage Identification	51
3.8.1. PtlGetUid	51

3.9. Process Identification	51
3.9.1. The Process Identification Type	52
3.9.2. PtlGetId	52
3.9.3. PtlGetPhysId	53
3.10. Memory Descriptors	53
3.10.1. The Memory Descriptor Type	54
3.10.2. The I/O Vector Type	55
3.10.3. PtlMDBind	56
3.10.4. PtlMDRelease	57
3.11. List Entries and Lists	57
3.11.1. The List Entry Type	59
3.11.2. PtlLEAppend	61
3.11.3. PtlLEUnlink	62
3.11.4. PtlLESearch	63
3.12. Match List Entries and Matching Lists	65
3.12.1. The Match List Entry Type	66
3.12.2. PtlMEAppend	70
3.12.3. PtlMEUnlink	71
3.12.4. PtlMESearch	72
3.13. Events and Event Queues	74
3.13.1. Kinds of Events	74
3.13.2. Event Occurrence	76
3.13.3. Failure Notification	77
3.13.4. The Event Structure	78
3.13.5. PtlEQAlloc	81
3.13.6. PtlEQFree	82
3.13.7. PtlEQGet	83
3.13.8. PtlEQWait	84
3.13.9. PtlEQPoll	84
3.14. Lightweight Counting Events	85
3.14.1. The Counting Event Type	86
3.14.2. PtlCTAlloc	86
3.14.3. PtlCTFree	87
3.14.4. PtlCTCancelTriggered	88

3.14.5. PtlCTGet	88
3.14.6. PtlCTWait	89
3.14.7. PtlCTPoll	90
3.14.8. PtlCTSet	91
3.14.9. PtlCTInc	91
3.15. Data Movement Operations	92
3.15.1. Portals Acknowledgment Type Definition	92
3.15.2. PtlPut	93
3.15.3. PtlGet	94
3.15.4. Portals Atomics Overview	95
3.15.5. PtlAtomic	98
3.15.6. PtlFetchAtomic	100
3.15.7. PtlSwap	102
3.15.8. PtlAtomicSync	103
3.16. Triggered Operations	104
3.16.1. PtlTriggeredPut	105
3.16.2. PtlTriggeredGet	106
3.16.3. PtlTriggeredAtomic	106
3.16.4. PtlTriggeredFetchAtomic	108
3.16.5. PtlTriggeredSwap	109
3.16.6. PtlTriggeredCTInc	110
3.16.7. PtlTriggeredCTSet	111
3.17. Deferred Communication Operations	111
3.17.1. PtlStartBundle	112
3.17.2. PtlEndBundle	113
3.18. Operations on Handles	113
3.18.1. PtlHandleIsEqual	113
3.19. Summary	114
4. Guide to Implementors	126
4.1. Run-time Support	126
4.2. Data Transfer	126
4.2.1. Sending Messages	126
4.2.2. Receiving Messages	129
4.3. Event Generation and Error Reporting	130

Bibliography	132
Appendices	134
A. Portals Design Guidelines	134
A.1. Mandatory Requirements	134
A.2. The <i>Will</i> Requirements	134
A.3. The <i>Should</i> Requirements	135
B. README Definition	136
C. Summary of Changes	137
C.1. Portals 4.3	137
C.2. Portals 4.2	138
C.3. Portals 4.1	139
C.4. Portals 4.0.2	140
C.5. Portals 4.0.1	141
C.6. Portals 4.0	141

List of Figures

- Figure 2-1. Graphical Conventions..... 21
- Figure 2-2. Portals Put (Send)..... 22
- Figure 2-3. Portals Get (Receive) from a match list entry 23
- Figure 2-4. Portals Get (Receive) from a list entry 24
- Figure 2-5. Portals Atomic Swap Operation 24
- Figure 2-6. Portals Atomic Sum Operation 25
- Figure 2-7. Simple Put Example 25
- Figure 2-8. Portals LE Addressing Structures 27
- Figure 2-9. Portals ME Addressing Structures 28
- Figure 2-10. Non-Matching Portals Address Translation..... 29
- Figure 2-11. Matching Portals Address Translation 31
- Figure 3-1. Portals Operations and Event Types 76

List of Tables

- Table 3-1. Object Type Codes 35
- Table 3-2. Event Type Summary 77
- Table 3-3. Event Field Definition 81
- Table 3-4. Legal Atomic Operation, Datatype, and Function Combinations 99
- Table 3-5. Portals Data Types 114
- Table 3-6. Portals Functions 116
- Table 3-7. Portals Return Codes 118
- Table 3-8. Portals Constants 118
- Table 4-1. Information Passed in a Send Request 127
- Table 4-2. Information Passed in an Acknowledgment 127
- Table 4-3. Information Passed in a “Counting” Acknowledgment 128
- Table 4-4. Information Passed in a Get Request 128
- Table 4-5. Information Passed in a Reply 129
- Table 4-6. Information Passed in an Atomic Request 129
- Table 4-7. Portals Operations and ME/LE Permission Flags 130

List of Implementation Notes

Note 1. No wire protocol.....	21
Note 2. Location of event queues and counting events.....	23
Note 3. Protected space.....	23
Note 4. Size of handle types.....	36
Note 5. Unique handles.....	37
Note 6. Memory descriptors that bind inaccessible memory.....	54
Note 7. Memory registration.....	57
Note 8. Optimization for Duplicate Memory Descriptors.....	57
Note 9. List entries that bind inaccessible memory.....	58
Note 10. PtILEUnlink() and unlinked handles.....	63
Note 11. Checking <i>match_id</i> Argument.....	71
Note 12. Completion of portals operations.....	78
Note 13. Size of event queue and reserved space.....	82
Note 14. Minimizing cost of counting events.....	86
Note 15. Portals Atomic Synchronization.....	103
Note 16. Ordering of Triggered Operations.....	104
Note 17. Purpose of Bundling.....	112

Preface

In the early 1990s, when memory-to-memory copying speeds were an order of magnitude faster than the maximum network bandwidth, it did not matter if data had to go through one or two intermediate buffers on its way from the network into user space. This began to change with early massively parallel processing (MPP) systems, such as the nCUBE-2 and the Intel Paragon, when network bandwidth became comparable to memory bandwidth. An intermediate memory-to-memory copy now meant that only half the available network bandwidth was used.

Early versions of Portals solved this problem in a novel way. Instead of waiting for data to arrive and then copy it into the final destination, Portals, in versions prior to 3.0, allowed a user to describe what should happen to incoming data by using data structures. A few basic data structures were used like Lego™ blocks to create more complex structures. The operating system kernel handling the data transfer read these structures when data began to arrive and determined where to place the incoming data. Users were allowed to create matching criteria and to specify precisely where data would eventually end up. The kernel, in turn, had the ability to DMA data directly into user space, which eliminated buffer space in kernel owned memory and slow memory-to-memory copies. We named that approach Portals Version 2.0. It was used until 2006 on the ASCI Red supercomputer, the first general-purpose machine to break the one teraflops barrier.

Although very successful on architectures with lightweight kernels, such as ASCI Red, Portals 2.0 proved difficult to port to Cplant [4] with its full-featured Linux kernel. Under Linux, memory was no longer physically contiguous in a one-to-one mapping with the kernel. This made it prohibitively expensive for the kernel to traverse data structures in user space. We wanted to keep the basic concept of using data structures to describe what should happen to incoming data. We put a thin application programming interface (API) over our data structures. We got rid of some never-used building blocks, improved some of the others, and Portals 3.0 was born [5].

Portals 3.0 evolved over three revisions to Portals 3.3 [19]. In the interim, the system context has changed significantly. Many newer systems are capable of offloading the vast majority of the Portals implementation to the network interface. Indeed, the rapid growth of bandwidth and available silicon area relative to the small decrease in memory latency has made it *desirable* to move latency sensitive tasks like Portals matching to dedicated hardware better suited to it. The implementation of Version 3.3 on ASC Red Storm (Cray XT3/XT4/XT5) illuminated many challenges that have arisen with these advances in technology. In this report, we document Portals 4 as a response to two specific challenges discovered on Red Storm. Foremost, while the performance of I/O buses has improved dramatically, the latency to cross an I/O bus has not fallen as dramatically as processor, memory and network performance has increased, negatively impacting target message rates. In addition, partitioned global address space (PGAS) models have risen in prominence and require lighter weight semantics compared to message passing.

Nomenclature

ACK	Acknowledgment.
FM	Illinois Fast Messages.
AM	Active Messages.
API	Application Programming Interface. A definition of the functions and semantics provided by library of functions.
ASCI	Advanced Simulation and Computing Initiative.
ASC	Advanced Simulation and Computing.
ASCI Red	Intel TeraFLOPS system installed at Sandia National Laboratories. First general-purpose system to break the one teraflops barrier.
CPU	Central Processing Unit.
DMA	Direct Memory Access.
EQ	Event Queue.
FIFO	First In, First Out.
FLOP	Floating Point Operation. (Also FLOPS or flops: Floating Point Operations per Second.)
GM	Glenn's Messages; Myricom's Myrinet API.
ID	Identifier.
Initiator	A <i>process</i> that initiates a message operation.
IOVEC	Input/Output Vector.
LE	List Entry.
MD	Memory Descriptor.
ME	Matching list Entry.
Message	An application-defined unit of data that is exchanged between <i>processes</i> .
Message Operation	Either a <i>put</i> operation, which writes data to a <i>target</i> , or a <i>get</i> operation, which reads data from a <i>target</i> , or an <i>atomic</i> operation, which updates data atomically.
MPI	Message Passing Interface.
MPP	Massively Parallel Processor.
NAL	Network Abstraction Layer.
NAND	Bitwise Not AND operation.
Network	A network provides point-to-point communication between <i>nodes</i> . Internally, a network may provide multiple routes between endpoints (to improve fault tolerance or to improve performance characteristics); however, multiple paths will not be exposed outside of the network.
NI	Abstract portals Network Interface.
NIC	Network Interface Card.
Node	A node is an endpoint in a <i>network</i> . Nodes provide processing capabilities and memory. A node may provide multiple processors (an SMP node). A node may also act as a <i>gateway</i> between networks.
OS	Operating System.
PM	Message passing layer for SCORED [12].
POSIX	Portable Operating System Interface.
Process	A context of execution. A process defines a virtual memory context. This context is not shared with other processes. Several threads may share the virtual memory context defined by a process.

RDMA	Remote Direct Memory Access.
RMPP	Reliable Message Passing Protocol.
SMP	Shared Memory Processor.
SUNMOS	Sandia national laboratories/University of New Mexico Operating System.
Target	A <i>process</i> that is acted upon by a message operation.
TCP/IP	Transmission Control Protocol/Internet Protocol.
Thread	A context of execution that shares a virtual memory context with other threads.
UDP	User Datagram Protocol.
UNIX	A multiuser, multitasking, portable OS.
VIA	Virtual Interface Architecture.

1. Introduction

1.1. Overview

This document describes the Portals network programming interface for communication between nodes in a system area network. Portals is designed to provide the building blocks necessary to create a diverse set of scalable, high performance application programming interfaces and language support run-times. The Portals API is designed to support a machine with millions of cores.

This document is divided into several sections:

Section 1 – Introduction.

The purpose and scope of the Portals API

Section 2 – An Overview of the Portals 4 API.

A brief overview of the Portals API, introducing the key concepts and terminology used in the description of the API

Section 3 – The Portals 4 API.

The functions and semantics of the Portals API in detail

Section 4 – Guide to Implementors.

A guide to implementors, highlighting subtleties of the standard that are critical to an implementation’s design

Appendix A – Portals Design Guidelines.

The guiding principles behind the Portals API design

Appendix B – README-template.

A template for a README file to be provided by each implementation

Appendix C – Summary of Changes.

A list of changes between versions since Portals 3.3

1.2. Purpose

Portals aims to provide a scalable, high performance network programming interface for High Performance Computing (HPC) systems. Portals provides an interface to support both the Message Passing Interface (MPI) [15] standard as well as the various partitioned global address space (PGAS) models, such as Unified Parallel C (UPC), Co-Array Fortran (CAF), and OpenSHMEM [10, 8]. While neither MPI nor PGAS models impose specific scalability limitations, many network programming interfaces do not provide the functionality needed to allow implementations of either model to reach scalability and performance goals.

The following are required properties of a network architecture to avoid scalability limitations:

- **Connectionless** – Many connection-oriented architectures, such as InfiniBand [11], VIA [9] and TCP/IP sockets, have practical limitations on the number of peer connections that can be established. In large-scale parallel systems, any node must be able to communicate with any other node without costly connection establishment and tear down.

- Network independence – Many communication systems depend on the host processor to perform operations in order for messages in the network to be consumed. Message consumption from the network should not be dependent on host processor activity, such as the operating system scheduler or user-level thread scheduler. Applications must be able to continue computing while data is moved in and out of the application’s memory.
- User-level flow control – Many communication systems manage flow control internally to avoid depleting resources, which can significantly impact performance as the number of communicating processes increases. While Portals provides building blocks to enable flow control (See Section 2.8), it is the responsibility of the application to manage flow control. An application should be able to provide final destination buffers into which the network can deposit data directly.
- OS bypass – High performance network communication should not involve memory copies into or out of a kernel-managed protocol stack. Because networks are now as fast as memory buses, data has to flow directly into user space.

The following are properties of a network architecture that avoid scalability limitations for an implementation of MPI:

- Receiver-managed data placement – Message passing implementations where the sender determines the target location of a data transmission require a persistent block of memory to be available for every process, requiring memory resources to increase with job size.
- User-level bypass (application bypass) – While OS bypass is necessary for high performance, it alone is not sufficient to support the *progress rule* of MPI asynchronous operations. After an application has posted a send or a receive, data must be delivered and acknowledged without further intervention from the application.
- Unexpected messages – Few communication systems have support for receiving messages for which there is no prior notification. Support for these types of messages is necessary to avoid flow control and protocol overhead.

1.3. Background

Portals was originally designed for and implemented on the nCUBE-2 machine as part of the SUNMOS (Sandia/UNM OS) [14] and Puma [20] lightweight kernel development projects. Portals went through three design phases [18], with the most recent one being used on the 13000-node (38,400 cores) Cray Red Storm [2] that became the Cray XT3/XT4/XT5 product line. Portals has been very successful in meeting the needs of such large machines, not only as a layer for a high-performance MPI implementation [7], but also for implementing the scalable run-time environment and parallel I/O capabilities of the machine.

The third-generation Portals implementation was designed for a system where the work required to process a message was long relative to the round trip between the application and the Portals data structures. However, in modern systems where processing is offloaded onto the network interface, the time to post a receive is dominated by the round trip across the I/O bus. This latency has become large relative to message latency and per message overheads (gap). This limitation was exposed by implementations on the Cray Red Storm system. Version 4.0 of Portals addresses this problem by adding the building blocks necessary to support the concept of *unexpected messages*. The second limitation exposed on Red Storm was the relative weight of handling newer PGAS programming models. PGAS programming models do not need the extensive matching semantics required by MPI and I/O libraries and can achieve significantly lower latency and higher message throughput without matching. Version 4.0 of Portals adds a lightweight, non-matching interface to support these semantics as well as lightweight events and acknowledgments. Finally, version 4.0 of Portals reduces the overheads in numerous implementation paths by simplifying events, reducing the size of acknowledgments, and generally specializing interfaces to eliminate functionality that experience has shown to be unnecessary. Version 4.3 is a refinement of Versions 4.0, 4.1 and 4.2, addressing issues found during the implementation of Portals 4.0, 4.1 and 4.2.

1.4. Scalability

The primary goal in the design of Portals is scalability. Portals is designed specifically for an implementation capable of supporting a parallel job running on millions of processing cores or more. Performance is critical only in terms of scalability. That is, the level of message passing performance is characterized by how far it allows an application to scale and not by how it performs in micro-benchmarks (e.g., a two-node bandwidth or latency test).

The Portals API is designed to allow for scalability, not to guarantee it. Portals cannot overcome the shortcomings of a poorly designed application program. Applications that have inherent scalability limitations, either through design or implementation, will not be transformed by Portals into scalable applications. Scalability must be addressed at all levels. Portals does not inhibit scalability and it does not guarantee it either. No Portals operation requires global communication or synchronization.

Similarly, a quality implementation is needed for Portals to be scalable. A non-scalable implementation, underlying network protocol, or hardware will result in a non-scalable Portals implementation and application.

To support scalability, the Portals interface maintains a minimal amount of state. By default, Portals provides reliable, ordered delivery of messages between pairs of processes. Portals is connectionless: a process is not required to explicitly establish a point-to-point connection with another process in order to communicate. Moreover, all buffers used in the transmission of messages are maintained in user space. The *target* process determines how to respond to incoming messages, and messages for which there are no buffers are discarded.

1.5. Communication Model

Portals combines the characteristics of both one-sided and two-sided communication. In addition to more traditional “put” and “get” operations, they define “matching put” and “matching get” operations. The destination of a *put* (or send) is not an explicit address; instead, messages target list entries (potentially with matching semantics or an offset) using the Portals addressing semantics that allow the receiver to determine where incoming messages should be placed. This flexibility allows Portals to support both traditional one-sided operations and two-sided send/receive operations.

Portals allows the *target* to determine whether incoming messages are acceptable. A *target* process can choose to accept message operations from a specific process or all processes, in addition to the ability to limit messages to a specified *initiator* usage id.

1.6. Zero Copy, OS Bypass, and Application Bypass

In traditional system architectures, network packets arrive at the network interface card (NIC), are passed through one or more protocol layers in the operating system, and are eventually copied into the address space of the application. As network bandwidth began to approach memory copy rates, reduction of memory copies became a critical concern. This concern led to the development of zero-copy message passing protocols in which message copies are eliminated or pipelined to avoid the loss of bandwidth.

A typical zero-copy protocol has the NIC generate an interrupt for the CPU when a message arrives from the network. The interrupt handler then controls the transfer of the incoming message into the address space of the appropriate application. The interrupt latency, the time from the initiation of an interrupt until the interrupt handler is running, is fairly significant. To avoid this cost, some modern NICs have processors that can be programmed to implement part of a message passing protocol. Given a properly designed protocol, it is possible to program the NIC to control the transfer of incoming messages without needing to interrupt the CPU. Because this strategy does not need to involve the OS on every message transfer, it is frequently called “OS bypass.” ST [21], VIA [9], FM [13], GM [17], PM [12], and Portals are examples of OS bypass mechanisms.

Many protocols that support OS bypass still require that the application actively participates in the protocol to ensure progress. As an example, the long message protocol of PM requires that the application receive and reply to a request to put or get a long message. This complicates the runtime environment, requiring a thread to process incoming requests, and significantly increases the latency required to initiate a long message protocol. Portals does not require activity on the part of the application to ensure progress. We use the term “application bypass” to refer to this aspect of Portals.

1.7. Faults

Reliable message transmission is challenging in modern high performance computing systems due to system scale, component failure rates, and application run-times. The Portals API recognizes that the underlying transport may not be able to successfully complete an operation once it has been initiated. This is reflected in the fact that the Portals API reports an event indicating the completion of every operation. Completion events indicate whether the operation completed successfully or not.

2. An Overview of the Portals API

In this chapter, we provide an overview of the Portals API and associated semantics. Detailed API functions and option definitions are presented in the next chapter.

2.1. Data Movement

A portal represents an opening in the address space of a process. Other processes can use a portal to read (*get*), write (*put*), or perform an atomic operation on the memory associated with the portal. Every data movement operation involves two processes, the *initiator* and the *target*. The *initiator* is the process that initiates the data movement operation. The *target* is the process that responds to the operation by accepting the data for a *put* operation, replying with the data for a *get* operation, or updating a memory location for, and potentially responding with the result from, an *atomic* operation.

In this discussion, activities attributed to a process may refer to activities that are actually performed by the process or *on behalf of the process*. The inclusiveness of our terminology is important in the context of *application bypass*. In particular, when we note that the *target* sends a reply in the case of a *get* operation, this is performed by Portals without the explicit involvement of the application. An implementation of Portals may use dedicated hardware, an operating system driver, a progress thread running in the application process, or some other option to generate the reply.

Figure 2-1 shows the graphical conventions used throughout this document. Some of the data structures created through the Portals API reside in user space to enhance scalability and performance, while others are kept in protected space for protection and to allow an implementation to place these structures into host or NIC memory. We use colors to distinguish between these elements.

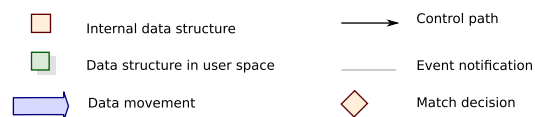


Figure 2-1. Graphical Conventions: Symbols, colors, and stylistic conventions used in the diagrams of this document.

Figures 2-2, 2-3, 2-4, and 2-5 present graphical interpretations of the Portals data movement operations: *put* (send), *get*, and *atomic* (the swap atomic is shown). In the case of a *put* operation, the *initiator* sends a put request ① message to the *target*. The *target* translates the portal addressing information in the request using its local portals structures. The data may be part of the same packet as the put request or it may be in separate packet(s) as shown in Figure 2-2. The Portals API does not specify a wire protocol. When the data ② has been put into the remote memory descriptor (or been discarded), the *target* optionally sends an acknowledgment ③ message.

**IMPLEMENTATION
NOTE 1:**

No wire protocol

This document does not specify a wire protocol. Portals requires a reliable communication layer with the semantics and progress rules specified in this document. Implementors are left great freedom in implementation design choices.

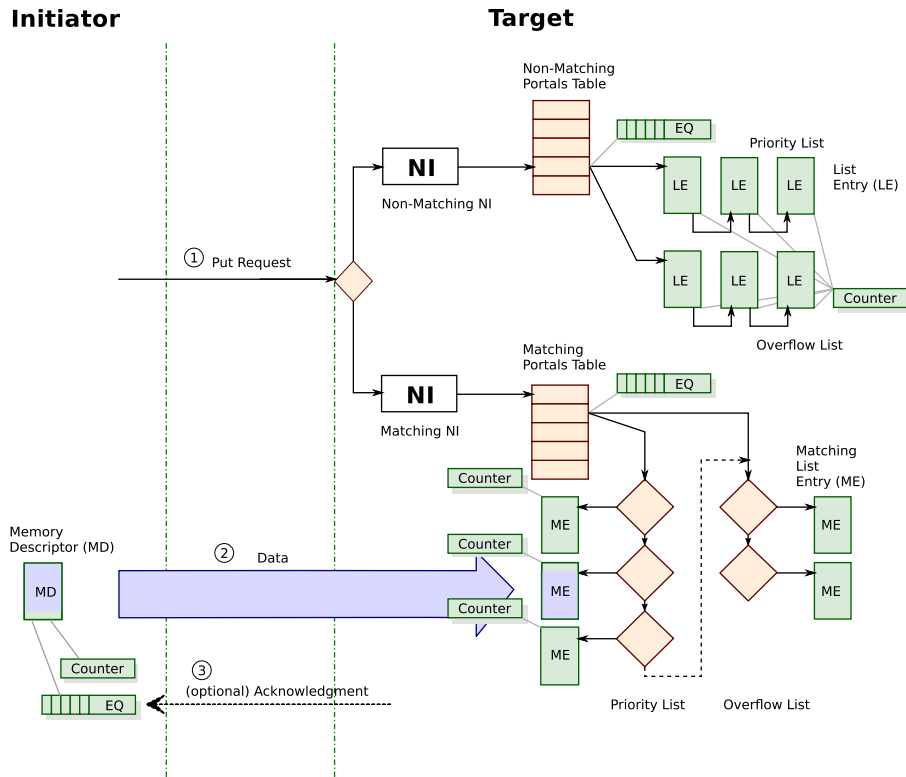


Figure 2-2. Portals Put (Send): Note that the put request ① is part of the header and the data ② is part of the body of a single message. Depending on the network hardware capabilities, the request and data may be sent in a single large packet or several smaller ones.

Figure 2-2 represents several important concepts in Portals 4. First, a message targets a *logical* network interface (NI) and a user may instantiate up to four logical network interfaces associated with a single *physical* network interface. A portals *physical* network interface is a per-process abstraction of a physical network interface (or group of interfaces). Logical network interfaces may be *matching* or *non-matching* and addressed by either *logical* (rank) or *physical* (nid/pid) identifiers. As indicated in Figure 2-2, separate logical network interfaces have independent resources. The second important concept illustrated in Figure 2-2 is that each portal table entry has three data structures attached: an event queue, a priority list, and an overflow list. The final concept illustrated in Figure 2-2 is that the overflow list is traversed after the priority list. If a message does not match in the priority list (matching interface) or it is empty (either interface), the overflow list is traversed.

Figure 2-2 illustrates another important Portals concept. The space the Portals data structures occupy is divided into protected and application (user) space, while the large data buffers reside in user space. Most of the Portals data structures reside in protected space. Often the Portals control structures reside inside the operating system kernel or the network interface card. However, they can also reside in a library or another process. See implementation note 2 for possible locations of the event queues.

IMPLEMENTATION NOTE 2: Location of event queues and counting events

Note that data structures that can only be accessed through the API, such as counting events and event queues, are intended to reside in user space. However, an implementation is free to place them anywhere it wants.

IMPLEMENTATION NOTE 3: Protected space

Protected space as shown for example in Figure 2-2 does not mean it has to reside inside the kernel or a different address space. The Portals implementation must guarantee that no alterations of Portals structures by the user can harm another process or the Portals implementation itself.

Figure 2-3 is a representation of a *get* operation from a *target* that does matching. The corresponding *get* from a non-matching *target* is shown in Figure 2-4. First, the *initiator* sends a request ① to the *target*. As with the *put* operation, the *target* translates the portals addressing information in the request using its local portals structures. Once it has translated the portals addressing information, the *target* sends a *reply* ② that includes the requested data.

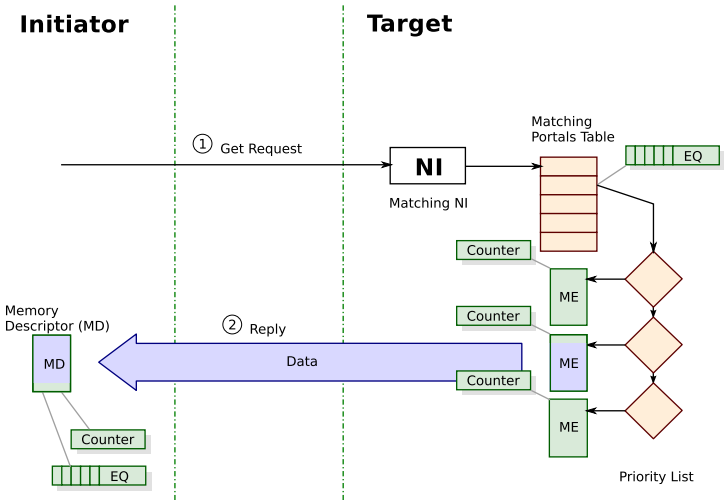


Figure 2-3. Portals Get from a match list entry.

Portals address translation (matching and permissions checks) is only performed at the *target* of an operation. Acknowledgments for *put* and *atomic* and replies to *get* and *atomic* operations bypass the portals address translation structures at the *initiator*. Acknowledgments and replies are only generated as the result of an action by the *initiator* and therefore do not require the same level of protection at the *initiator* as would be required at the *target*.

The third operation type, *atomic*, is depicted in Figure 2-5 for the swap operation and Figure 2-6 for a summation.

For the swap operation shown in Figure 2-5, the *initiator* sends a request ①, containing the *put* data and the operand value ②, to the *target*. The *target* traverses the local portals structures based on the information in the request to find the appropriate user buffer. The *target* then sends the *get* data in a *reply* message ③ back to the *initiator* and deposits the *put* data in the user buffer.

The sum operation shown in Figure 2-6 adds the *put* data into the memory region described by the list entry. The figure shows an optional *acknowledgment* sent back. The target data prior to the summation is not sent back, since the *initiator* used `PtlAtomic()` instead of `PtlFetchAtomic()`.

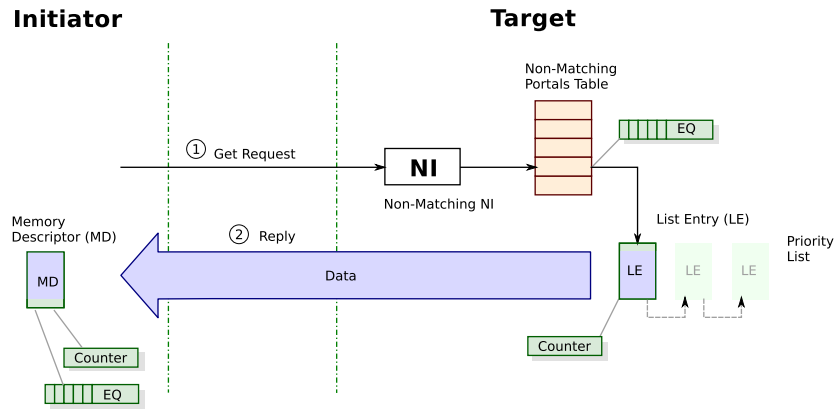


Figure 2-4. Portals Get from a list entry. Note that the first LE will be selected to reply to the *get* request.

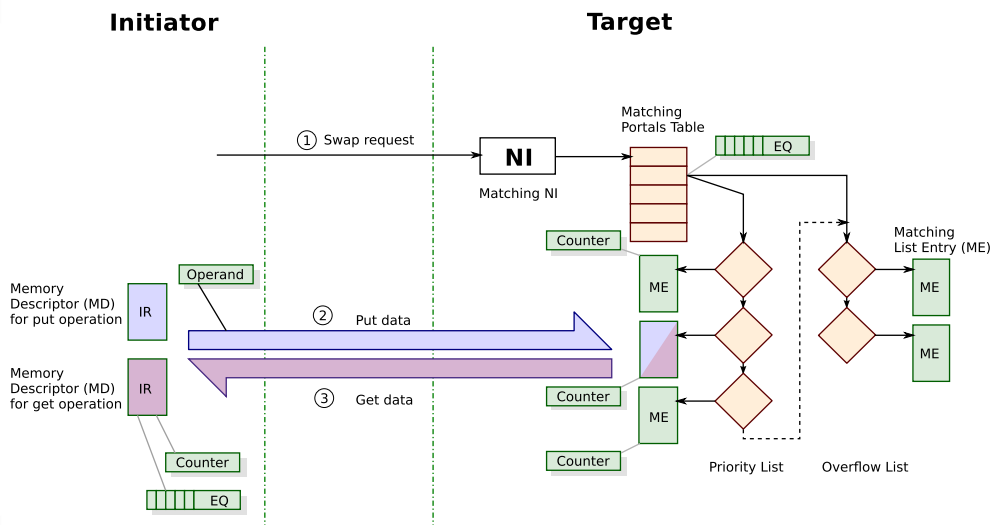


Figure 2-5. Portals Atomic (swap is shown). An atomic swap in memory described by a match list entry using an initiator-side operand.

2.2. Unreliable Datagrams

In certain cases, reliable delivery of network data is not required. To save resources and avoid network contention caused by retransmissions, an unreliable datagram mode is available. If the memory descriptor `PTL_MD_UNRELIABLE` option is set, guaranteed delivery of messages is not provided and messages may be dropped by the network. Only the *put* operation with no acknowledgment is supported in this mode.

2.3. Usage

Some of the diagrams presented in this chapter may seem daunting at first sight. However, many of the diagrams show all possible options and features of the Portals building blocks. In actual use, only some of them are needed to

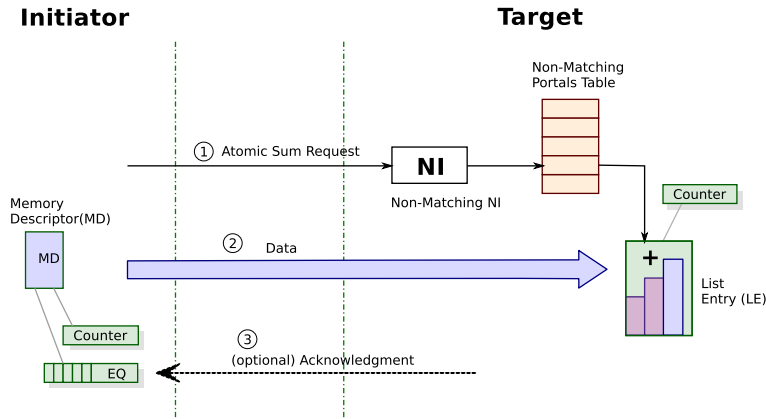


Figure 2-6. Portals Atomic (sum is shown). An atomic sum operation in memory described by a list entry.

accomplish a given function. Rarely will they all be active and used at the same time.

Figure 2-2 shows the complete set of options available for a *put* operation. In practice, a diagram like Figure 2-7 is much more realistic. It shows the Portals structures used to setup a one-sided *put* operation. A user of Portals needs to specify an initiator region where the data is to be taken from, and an unmatched target region to put the data. Offsets can be used to address portions of each region; e.g., a word at a time, and an event queue or a counting event inform the user when an individual transfer has completed.

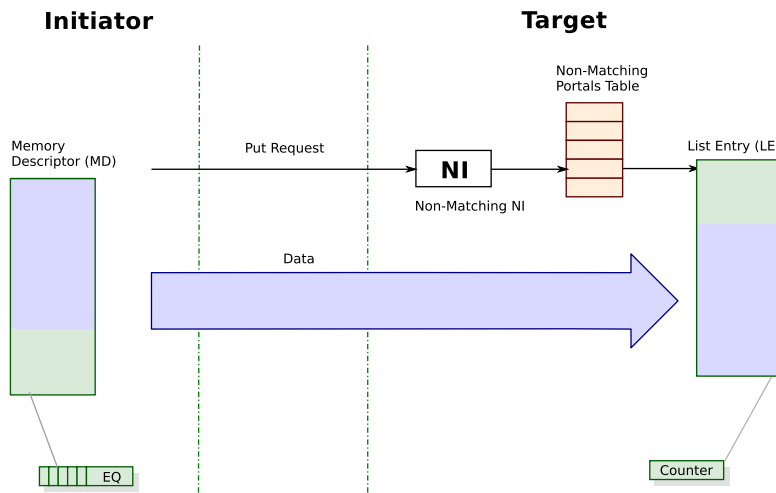


Figure 2-7. Simple Put Example: Not every option or Portals feature is needed to accomplish simple tasks such as the transfer of data from an initiator region to a target region.

Another example is Figure 2-6 which is simpler than Figure 2-5 and probably more likely to be used in practice. Atomic operations, such as the one in Figure 2-6 are much more likely to use a single unmatched target region. Such simple constructs can be used to implement global reference counters, or access locks.

2.4. Completion Events

Portals provides two mechanisms for recording completion events: full events (Section 3.13) and counting events (Section 3.14). Full events provide a complete picture of the transaction, including what type of event occurred, which buffer was manipulated, and identifying any errors that occurred. The full event can also carry a small amount of local data and, on the target, a small amount of out-of-band header data. Counting events, on the other hand, are designed to be lightweight and provide only a count of successful and failed operations (or successful bytes delivered). The delivery of events (full events or counting events) may be manipulated when creating a number of other structures.

2.5. Portals Addressing

One-sided data movement models (e.g., OpenSHMEM [8], SHMEM [10], ST [21], and MPI RMA [16]) typically use a process identifier and remote address to identify a memory address on a remote node. In some cases, the remote address is specified as a memory buffer identifier and offset. The process identifier identifies the *target* process, the memory buffer identifier specifies the region of memory to be used for the operation, and the offset specifies an offset within the memory buffer.

Portals lists provide one-sided addressing capabilities. Portals list entries serve as a memory buffer identifier that may be persistent or optionally removed from the list after a single use. Traditional one-sided addressing capabilities have proven to be a poor fit for tagged messaging interfaces, such as the Message Passing Interface [6]. To overcome these limitations, Portals also supports match list entries, which include additional semantics for receiver-managed data placement. Matching semantics are discussed in Section 2.5.2.

In addition to matching a pre-posted list entry, an incoming message also must pass a permissions check. The permissions check is *not* a component of identifying the correct buffer. It is *only* applied after the correct buffer has been identified. The permissions check has two components: the target of the message must allow the initiator to access the buffer and must allow the specified operation type. Each list entry and match list entry specifies which types of operations are allowed—put and/or get—as well as a usage ID that can be used to identify which initiators are allowed to access the buffer. A failure of the permissions check for an incoming message does not modify the Portals state in any way, except to update the status registers (see Section 3.3.7), and the message itself is discarded. Permissions IDs such as a usage ID must be contained in a protected header in the portals message. A protected header is part of a portals message where the user may not modify data, such that the usage ID inserted in the header is the same as the one allocated by the implementation.

Figures 2-8 and 2-9 are graphical representations of the structures used by a *target* in the interpretation of a portals address. The initiator's physical network interface and the specified target node identifier are used to route the message to the appropriate node and physical network interface. This logic is not reflected in the diagrams. The *initiator*'s logical network interface and the specified target process ID¹ are used to select the correct *target* process and the logical network interface. Each logical network interface includes a single portal table used to direct message delivery.

Discussion: *Portals loosely defines the concept of a physical network interface. A physical network interface may be a single hardware network interface or it may represent a collection of hardware network interfaces, with multi-rail support implemented within the Portals implementation.*

For example, in a system like BlueGene/L [1], an implementation may expose a physical network interface for the high speed network and another physical network interface for the Ethernet support and I/O network. On the other hand, a system with multiple InfiniBand HCAs may choose to expose a single physical network interface which load balances between the hardware interfaces. In both cases, a portal table will be created for each initialized logical network interface over each physical network interface for each process.

¹A logical *rank* can be substituted for the combination of node ID and process ID when logical endpoint addressing is used.

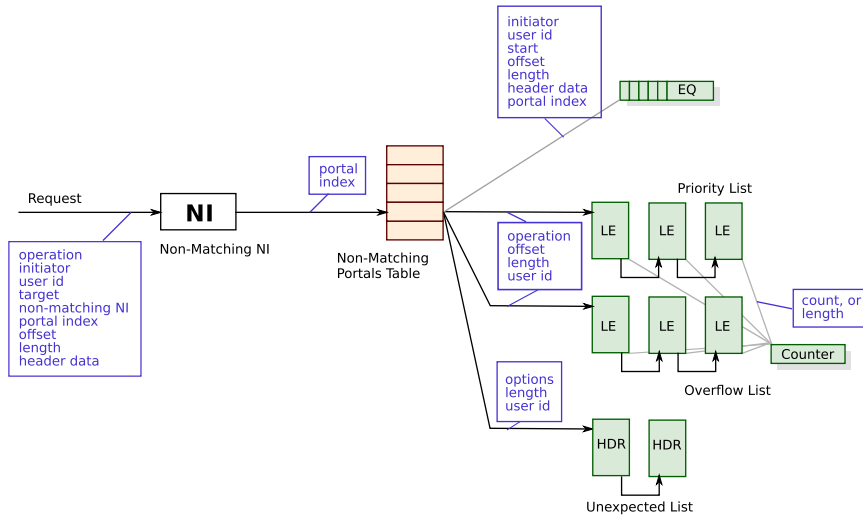


Figure 2-8. Portals Non-Matching Addressing Structures: The example shows the flow of information for a non-matched request at a target. Various pieces of information from the incoming header flow to the Portals structures where they are needed to process the request.

An initiator-specified portal index is used to select an entry in the portal table. Each entry of the portal table identifies three lists and, optionally, an event queue. The priority list and overflow list provide lists of remotely accessible address regions. Applications may append new list entries to either list, allowing complex delivery mechanisms to be built. Incoming messages are first processed according to the priority list and, if no matching entry was found in the priority list, are then processed according to the overflow list. In addition to providing an insertion point in the middle of the combined list structures by allowing insertions at the end of both the priority and overflow lists, the overflow list carries additional semantics to allow unexpected message processing.

The third list that is associated with each portal index is more transparent to the user and provides the building blocks for supporting unexpected messages. Each time a message is delivered into the overflow list, its header is linked into the unexpected list. The user cannot insert a header into the unexpected list, but can search the list for matching entries and, optionally, delete the matching entries from the list. Further, when a new list entry is appended to the priority list, the unexpected list is first searched for a match. If a match is found (i.e., had the list entry been on the priority list when the message arrived, the message would have been delivered into that list entry), the list entry is not inserted, the header is removed from the unexpected list (see Section 3.12 for details regarding exceptions to this behavior), and the application is notified a match was found in the unexpected list. A list entry in the overflow list may disable the use of the unexpected list for messages delivered into that list entry. All unexpected messages associated with a list entry must be handled by posting matching list entries in the priority list or searching and deleting prior to **PtlLEUnlink()** or **PtlMEUnlink()** successfully unlinking the overflow list entry. Unlike incoming messages, no permissions check is performed during the search of the unexpected queue. Therefore, the user is responsible for ensuring that the overflow list provides sufficient protection to memory and any further permissions checks must be performed by the user based on the overflow event data.

Each data manipulation event (e.g., **PTL_EVENT_PUT**) has a corresponding overflow event (e.g., **PTL_EVENT_PUT_OVERFLOW**) which is generated when a matching header is found in the unexpected list during list entry insertion. Overflow events may only occur after the data has been fully delivered to the overflow buffer. The overflow full event includes sufficient information (event type, start address, length, etc.) to determine what operation occurred and where the data was delivered into the overflow list. If the *mlength* in the full event is less than the *rlength*, the message was truncated. It is the responsibility of the application to retrieve the message body, if necessary. For cases where an application posts a **PtlMEAppend()** or **PtlLEAppend()** that does not provide a large enough buffer for the match entry in the overflow list, the application must check the returned *mlength* against the size of the posted buffer to ensure that truncation did not happen. Truncation checks only occur on the overflow list

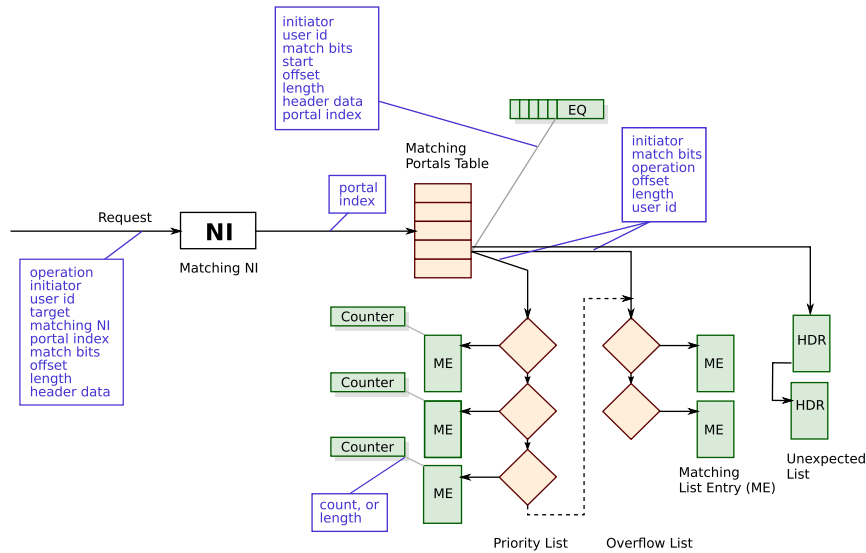


Figure 2-9. Portals Matching Addressing Structures: The example shows the flow of information for a matched request at a target. Various pieces of information from the incoming header flow to the Portals structures where they are needed to process the request.

entry when it is matched, they are not re-performed on checks of posted operations that match in the overflow list.

If the incoming message is not delivered into either the priority or overflow list and flow control is not enabled on the portal table entry, the message is discarded and the `PTL_SR_DROP_COUNT` status register is incremented (see Section 3.3.7). If flow control is enabled on the portal table entry, flow control is triggered and a `PTL_EVENT_PT_DISABLED` full event is generated in the event queue associated with the portal table entry (see Section 2.8).

In typical scenarios, MPI point-to-point communication uses the matching interface and full events, while OpenSHMEM uses the non-matching interface and lightweight counting events. The overflow list may act as either a building block for handling MPI unexpected messages (when the unexpected list is enabled) or as a mechanism for allowing insertion into the middle of a list (when the unexpected list is disabled).

2.5.1. Lists and List Entries

Lists and list entries provide semantics similar to that found in traditional one-sided interfaces. List entries identify a memory region as well as an optional counting event. The memory region specifies the memory to be used in the operation, and the counting event is optionally used to record the occurrence of operations. Information about the operations is (optionally) recorded in the event queue attached to the portal table entry.

Figure 2-10 shows the logical flow of Portals address translation on a non-matching logical network interface. The first list entry (LE) in a list *always* matches. Authentication is provided through fields associated with the LE and act as *permission* fields, which can cause the operation to fail. An operation can fail to fit in the region provided and, if so, will be truncated. Other semantics provided by match list entries—such as locally managed offsets—are not supported. The overflow list is checked after the priority list, if necessary. The non-matching translation path has the same event semantics as a matching interface. The important difference between the non-matching interface and the matching interface is that the Portals address translation semantics for the non-matching interface always match the first entry. This allows fully pipelined operation for the non-matching address translation.

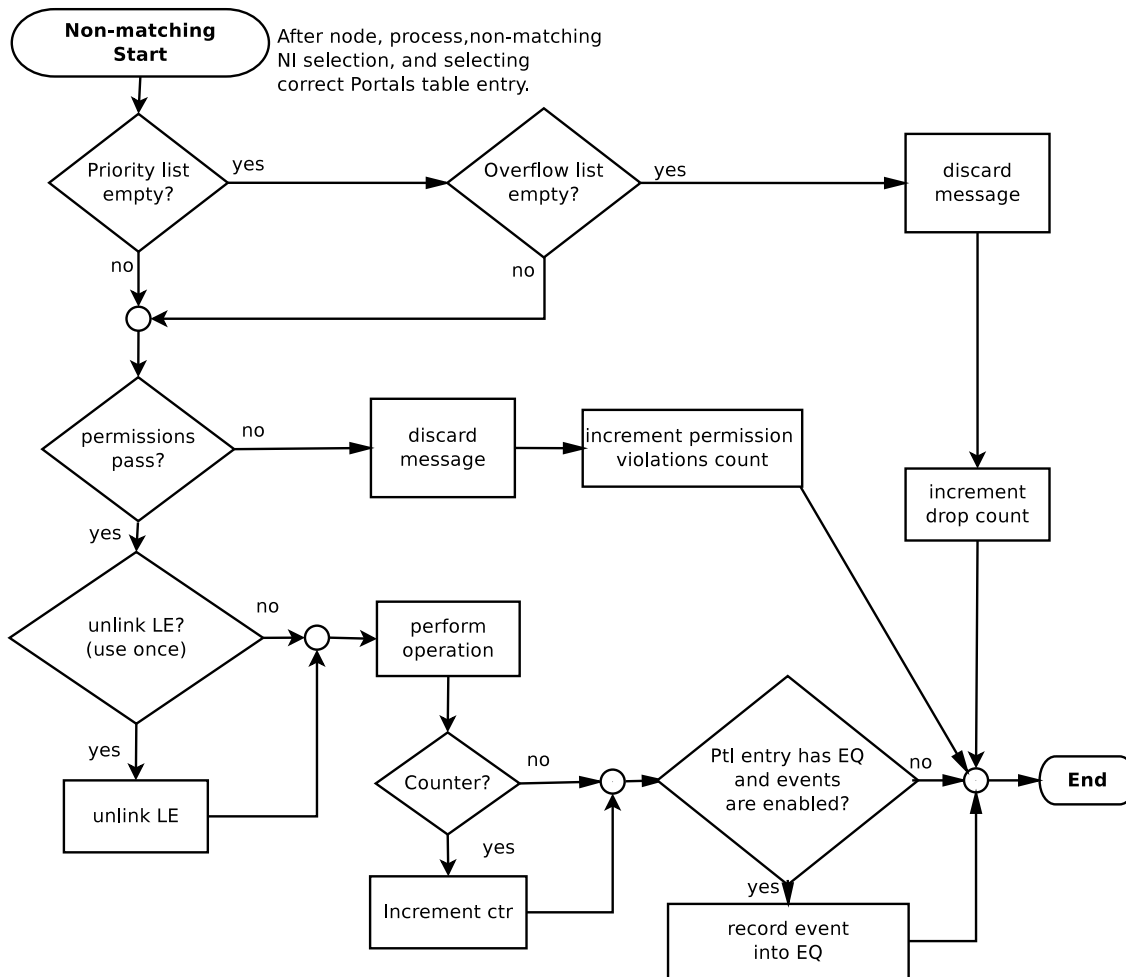


Figure 2-10. Non-Matching Portals Address Translation.

Discussion: List entries may be persistent or automatically unlink after first use. Implementations may be able to provide much higher message rates if the priority list contains a persistent list entry at the head of the list. One-sided programming interfaces such as OpenSHMEM and MPI-3 one-sided should be able to take advantage of this performance gain.

2.5.2. Match Lists and Match List Entries

In addition to the standard address components (process identifier, memory buffer identifier, and offset), a portals address can include information identifying the *initiator* (source) of the message and a set of match bits. This addressing model is appropriate for supporting traditional two-sided message passing operations. Specifically, the Portals API provides the flexibility needed for an efficient implementation of MPI-1, which defines two-sided operations, with one-sided completion semantics.

For a matching logical network interface, each match list entry specifies two bit patterns: a set of “do not care” bits (ignore bits) and a set of “must match” bits (match bits). Along with the source node ID (NID) and the source process ID (PID), these bits are used in a matching function to select the correct match list entry. In addition, if truncation is disabled (PTL_ME_NO_TRUNCATE is set), the message must fit in the buffer. If the message does not fit, the message does not match that entry and matching continues with the next entry.

In addition to *initiator*-specified offsets, match list entries also support locally managed offsets, which allow efficient packing of multiple messages into a single match list entry. When locally managed offsets are enabled, the initiator-specified offset is ignored. A match list entry may additionally specify a minimum available space threshold (*min_free*), after which a persistent match list entry is automatically unlinked. The combination of locally managed offsets, minimum free thresholds, and overflow list semantics facilitate the efficient implementation of MPI unexpected messages.

Figure 2-11 illustrates the steps involved in translating a portals address when matching is enabled, starting from the first element in a priority list. If the match criteria specified in the match list entry are met, the permissions check passes, and the match list entry accepts the operation, the operation (*put*, *get*, or *atomic*) is performed using the memory region specified in the match list entry. Note that matching is done using the match bits, ignore bits, and either the node identifier and process identifier or the rank.

If the match list entry specifies that it is to be unlinked based on the *min_free* semantic or if it is a use once match list entry, the match list entry is removed from the match list, and the resources associated with the match list entry are reclaimed. If there is an event queue specified in the portal table entry and the match list entry accepts the full event, the operation is logged in the event queue. An event is delivered when no more actions, as part of the current operation, will be performed on this match list entry.

If the match criteria specified in the match list entry are not met, the address translation continues with the next match list entry. If the end of the priority list has been reached, address translation continues with the overflow list. Once a matching match list entry has been identified, if the permissions check fails or the match list entry rejects the operation, the matching ceases and the message is dropped without modifying the list state.

2.6. Modifying Data Buffers

Users pass data buffers into the Portals implementation as either a source of data or the destination of data. For buffers where data is being delivered (e.g., at the *target*, or in a reply buffer at the *initiator*), the Portals API allows user memory to be used as a scratch space as long as the operation is larger than *max_atomic_size*. That means an implementation can utilize user memory as scratch space and staging buffers for operations larger than this threshold. When the operation is larger than *max_atomic_size*, the user memory is not guaranteed to reflect exactly the data that has arrived until the operation succeeds and the event is delivered. In fact, for operations larger than *max_atomic_size*, the memory may be changed in unpredictable ways while the operation is progressing. Once the operation completes, the memory associated with the operation will not be subject to further modification (from this operation). Notice that unsuccessful operations may alter memory used to receive data in an essentially unpredictable fashion.

The Portals API explicitly prohibits modifying the buffer passed into a *put*. Similarly, an implementation must not alter data in a user buffer that is used in a *reply* until the operation completes. This is independent of whether the operation succeeds or fails.

2.7. Ordering

There are three types of ordering typically defined by higher-level languages and message passing APIs: message ordering, data ordering, and write ordering. The message ordering definition controls the order in which messages are processed by the match engine between a pair of endpoints. The data ordering definition controls the order in which the data of two different messages is delivered into memory. The write ordering definition controls the order in which the bytes of a single message are written to memory. Ordering is a complex subject with a variety of high-level definitions in programming languages and message passing APIs. Portals does not define any write ordering, but it has a variety of options to control message and data ordering. As a general overview, Portals guarantees byte-granularity data ordering for short messages between a pair of endpoints when targeting a specific list entry or match list entry. For all messages regardless of size, message ordering is provided unless it is disabled using the

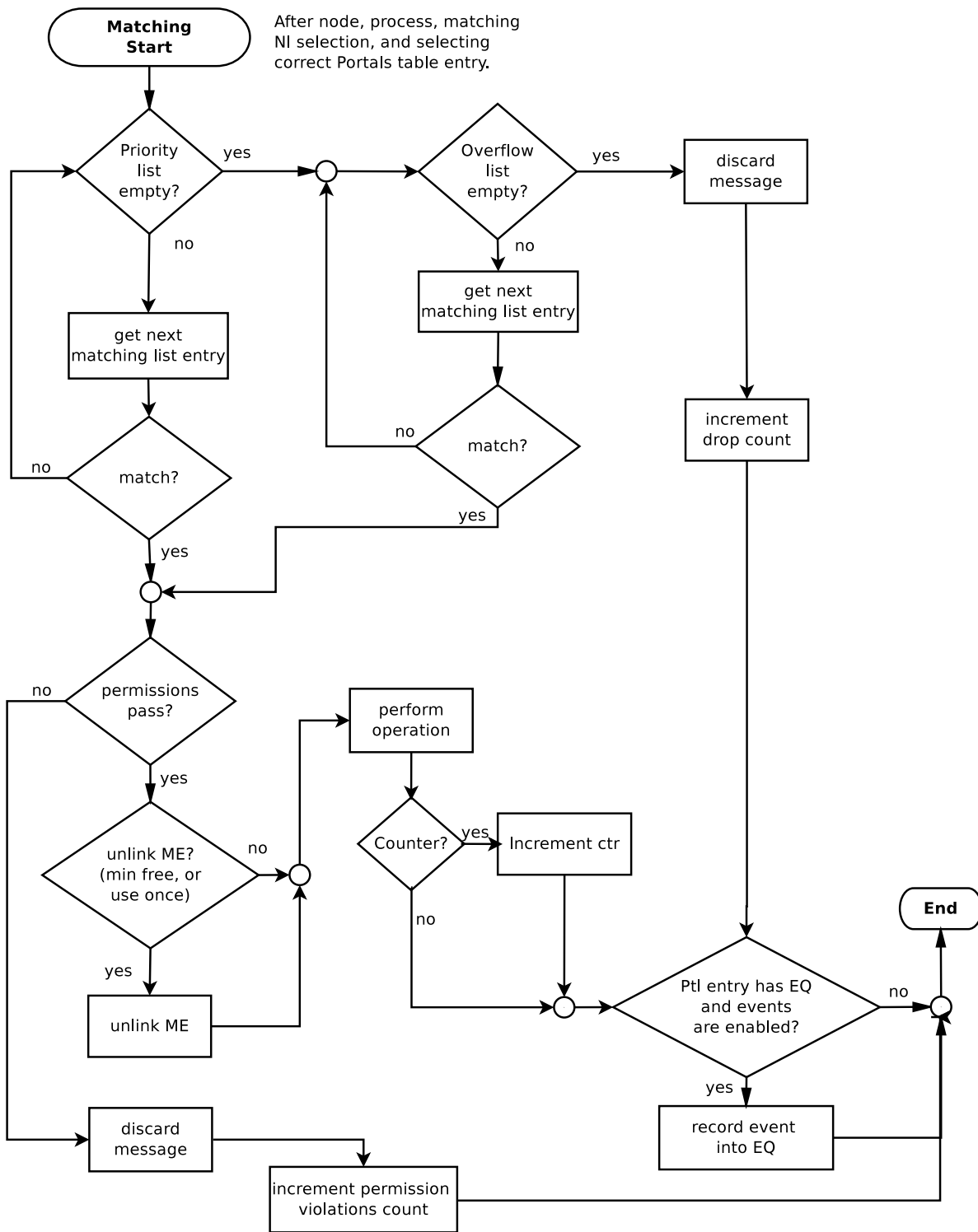


Figure 2-11. Matching Portals Address Translation.

PTL_MD_UNORDERED option in the `ptl_md_t`. This supports the MPI two-sided message ordering requirements while providing the flexibility to disable ordering when it is not needed.

2.7.1. Short Message Ordering Semantics

The default ordering semantics for Portals messages differ for short and long messages. The threshold between “short” and “long” is defined by two parameters, the maximum write-after-write and read-after-write sizes (*max_waw_ordered_size*) and the maximum write-after-read size (*max_war_ordered_size*). Both parameters are controlled by the *desired* and *actual* arguments of `PtINIInit()`. Note that replies and acknowledgments do not require ordering.

When one message that stores data (*put, atomic*) is followed by a message that stores data or retrieves data (*put, atomic, get*) from the same initiator to the same target and both messages are less than the *max_waw_ordered_size* in length, a byte from the second message that targets the same offset within the *same* LE (or ME) as a byte from the first message will perform its access after the byte from the first message. Similarly, when one message that retrieves data (*get*) is followed by a second message that stores data (*put, atomic*) from the same initiator to the same target and both messages are less than *max_war_ordered_size* in length, a byte from the second message that targets the same offset within the same LE (or ME) as a byte from the first message will perform its access after the byte from the first message.

Even for small messages, Portals does not guarantee write ordering, i.e., the order in which individual bytes of a single message are delivered is always unspecified. In addition, the order in which non-overlapping bytes of two different messages are written is *not* specified unless the implementation provides total data ordering and both target and initiator NIs have the `PTL_TOTAL_DATA_ORDERING` option set in the *actual_features* limits field. If either the initiator or target do not have `PTL_TOTAL_DATA_ORDERING` set, total data ordering will not be provided. When total data ordering is provided and the short message constraints are met, the first message must be entirely delivered before any part of the second message is delivered. Total data ordering provides additional ordering guarantees over *max_waw_ordered_size* in that it requires that the entire first message must be written in its entirety before the second message can be written (including non overlapping bytes), while write-after-write ordering only guarantees ordering within overlapping offset regions within the same LE or ME. Support for the ordering of bytes between messages is an optional feature, since some implementations may be unable to provide such strict ordering semantics. It is the responsibility of the initiator to ensure that the target provides total data ordering. Targets are free to not provide total data ordering support for messages incoming from `PTL_TOTAL_DATA_ORDERING` enabled initiator NIs.

2.7.2. Long Message Ordering Semantics

The default ordering semantics for Portals messages that have a length that is longer than the *max_waw_ordered_size* (or *max_war_ordered_size*, as appropriate) are much weaker. For long messages, the ordering semantics only require that messages sent between a pair of processes are matched at the target in the order they were sent. The underlying implementation is free to deliver the *body* of two messages in whatever order is necessary. This provides additional flexibility to the underlying implementation. For example, the implementation can use a retransmission protocol that only retransmits a portion of a lost message without violating ordering. Similarly, an implementation is free to use adaptive routing to deliver the body of the message. Note that replies and acknowledgments do not require ordering.

Discussion: *The specified ordering semantics of Portals are not necessarily sufficient to allow a `shmem_fence()` operation to be treated as a no-op. Portals only guarantees ordering semantics sufficient for `shmem_fence()` to be a no-op when `PTL_TOTAL_DATA_ORDERING` is returned in the options field of the *actual* limits and the operations are both shorter than *max_waw_ordered_size*.*

2.7.3. Relative Ordering of Operations in Overlapping Portals

The result of a *put* or *atomic* operation transferring data from a memory location (within a memory descriptor) which is currently the target of a remote operation (within a list entry) is undefined. Data is only available for transmit after an event indicating the completion of the arriving message has been delivered. Triggered operations are ordered

safely, since they do not trigger until the counting event is delivered and are triggered in the order in which they were posted for a given counter. Operations simultaneously taking place on overlapping portals have the same behavior as overlapping operations on a single portal (e.g. LE) without `PTL_TOTAL_DATA_ORDERING`; see Section 2.7.1 for more information on overlapping operations on an LE.

2.7.4. Ordering of Unexpected Messages

Messages delivered into an overflow list entry have the same ordering semantics as messages delivered into priority list entries. Data delivery into a overflow list entry is ordered like that of a priority list entry. Overflow list entries are typically programmed to capture the headers (and possibly buffer the payload) of messages that arrived before a matching entry could be appended to the priority list. Information on the ordering of unexpected headers is captured to enable MPI matching ordering semantics.

2.7.5. Relaxing Message Ordering

In many modern networks, adaptive routing can be used to improve the overall network throughput. For these networks, it may be useful for the application to express to the implementation when it is possible to relax the ordering on messages. Portals provides two mechanisms to relax ordering. First, when the application calls `PtINIInit()`, it can request a `max_waw_ordered_size` and `max_war_ordered_size` of zero in `ptl_ni_limits_t` (see Sections 3.6.1 and 3.6.2). This informs the implementation that data ordering is not needed (e.g., in the two sided semantics for MPI). Second, the application can set the `PTL_MD_UNORDERED` option on the `ptl_md_t` used to send the data (see Section 3.10). This turns off both message and data ordering.

2.8. Flow Control

Historically, on some large machines, MPI over Portals has run into problems where the number of unexpected messages has caused the exhaustion of event queue space or buffer space set aside for unexpected messages. MPI implementations over past versions of Portals have handled the overflow by aborting the application. Other networks use “receiver not ready” NACKs and retransmits at the hardware level. Unfortunately, this is known to impact pipelining in the NIC. In attempting to address this challenge, Portals 4 adopts the philosophy that resource exhaustion is an exceptional operating mode and recovery may be slow, but must be possible.

When resources are exhausted, whether they are user allocated resources like EQ entries or implementation level resources, the implementation may choose to block new message processing for a constrained amount of time. If the resources remain exhausted, the behavior of Portals depends on the type of operation which caused the exhaustion and, potentially, options set by the user.

A local operation which generates events (such as a call to `PtILEAppend()` or `PtIMEAppend()`) or a response from a target-side operation (such as an acknowledgment or reply) is not required to trigger flow control and may cause the event queue to overflow, resulting in dropped events. An implementation may choose to trigger flow control for local operations, but is not required to do so.

Discussion: *The user must use some care when posting a new list entry to ensure that local events do not overflow the event queue. A sufficiently large event queue, drained before posting the list entry, will provide sufficient protection. Implementations may choose to perform more resource exhaustion checking to prevent overflowing the event queue, but are not required to do so.*

A target-side operation (such as the processing of an incoming `put` or `get` operation) which targets a portal table entry on which the `PTL_PT_FLOWCTRL` option has not been set will not trigger flow control. If the message failed to match in the priority or overflow lists or the message would have matched in the overflow list and the unexpected headers list is full, the message will be dropped, with the `PTL_SR_DROP_COUNT` status register incremented as specified in

Section 2.5. An acknowledgment or reply event will not be generated in this case. If the constrained resource was an event queue, the message will be delivered and any acknowledgment or reply will be generated, but the target-side event will be lost.

A target-side operation (such as the processing of an incoming *put* or *get* operation) which targets a portal table entry on which the `PTL_PT_FLOWCTRL` option has been set will trigger flow control. When flow control is triggered, the implementation must disable the portal table entry and deliver a `PTL_EVENT_PT_DISABLED` full event to the application (See Implementation Note 13). At this point, all messages targeting that portal table entry for that process must be dropped until `PtIPTEnable()` is called, including the message that caused the flow control event. Messages that are dropped due to a flow control event do not modify any portion of the buffer described by the target list entry or match list entry. In addition, the `PTL_EVENT_ACK` or `PTL_EVENT_REPLY` event associated with that message (and subsequent in flight messages) indicate failure. The *ni_fail_type* of any generated full event must be `PTL_NI_PT_DISABLED`. While a disabled portal table entry will refuse any communication operations, local operations (such as `PtILEAppend()`) continue to be processed.

Discussion: *It is important to note that remote flow control failure notification is only delivered to the initiator of an operation in the `PTL_EVENT_ACK` or `PTL_EVENT_REPLY` event; thus, it is necessary for a user to request acknowledgments at the initiator to be notified of a flow control situation.*

While any internal, potentially implementation specific, resource exhaustion can cause a flow control event, three Portals level resource exhaustion types must cause a flow control event when they occur. If flow control is enabled, the following three scenarios must invoke flow control. First, if an event queue attached to a portal table entry is full and the message would generate a full event, flow control must be invoked. Second, if a message arrives at a portal table entry and does not find a match in either the priority list or the overflow list, flow control must be invoked. Finally, if the space available to buffer unexpected message headers is exhausted (e.g., as indicated by *max_unexpected_headers*), flow control must be invoked.

Discussion: *The application must be involved in flow control recovery. The difficulty in recovering is largely driven by the ordering constraints of the application. Interfaces with loose ordering semantics (such as GASNet) may be able to reduce resource utilization and re-enable a portal table entry without any global communication. Strictly ordered interfaces, such as MPI, must quiesce the library, ensure that resources are available, reach a global consensus that the network is quiesced (likely using another portal table entry for communication), re-enable the portal table entry, and restart communication. Quiescing the library requires the MPI library to insure that no more messages are in flight targeting the node that has experienced resource exhaustion. Making resources available involves draining all full events from the event queue associated with the portal table entry, replenishing the user allocated buffers on the overflow list, and draining unexpected messages from the Portals implementation.*

2.9. Multi-Threaded Applications

The Portals API supports a generic view of multi-threaded applications. From the perspective of the Portals API, an application program is defined by a set of processes. Each process defines a unique address space. The Portals API defines access to this address space from other processes (using portals addressing and the data movement operations). A process may have one or more *threads* executing in its address space.

With the exception of waiting (`PtIEQWait()`, `PtICTWait()`), polling (`PtIEQPoll()`, `PtICTPoll()`), portal table manipulation functions (`PtIPTDisable()`, `PtIPTEnable()`), and some allocation routines (such as `PtICTAlloc()`, `PtICTFree()`, `PtIEQAlloc()`, `PtIEQFree()`, `PtIMEUnlink()`), every function in the portals API is non-blocking. Every function in the Portals API is atomic with respect to both other threads and external operations that result from data movement operations. While individual operations are atomic, sequences of these operations may be interleaved between different threads and with external operations. In other words, calls into the Portals API are thread safe. The Portals API does not provide any mechanisms to control this interleaving. It is expected that these mechanisms will be provided by the API used to create threads.

3. The Portals API

3.1. Naming Conventions and Typeface Usage

The Portals API defines four types of entities: functions, types, return codes, and constants. Functions always start with **Ptl** and use mixed upper and lower case. When used in the body of this report, function names appear in sans serif bold face, e.g., **PtlInit()**. The functions associated with an object type will have names that start with **Ptl**, followed by the two letter object type code shown in column *yy* in Table 3-1. As an example, the function **PtlEQAlloc()** allocates resources for an event queue.

Table 3-1. Object Type Codes.

<i>yy</i>	<i>xx</i>	Name	Section
NI	ni	Network Interface	3.6
PT	pt	Portal Table Entry	3.7
MD	md	Memory Descriptor	3.10
LE	le	List Entry	3.11
ME	me	Matching list Entry	3.12
EQ	eq	Event Queue	3.13
CT	ct	Count	3.14

Type names use lower case with underscores to separate words. Each type name starts with **ptl_** and ends with **_t**. When used in the body of this report, type names appear like this: **ptl_match_bits_t**.

Return codes start with the characters **PTL_** and appear like this: **PTL_OK**.

Names for constants use upper case with underscores to separate words. Each constant name starts with **PTL_**. When used in the body of this report, constant names appear like this: **PTL_ACK_REQ**.

The definition of named constants, function prototypes, and type definitions must be supplied in a file named `portals4.h` that can be included by programs using Portals. Implementations must also provide the same interface in a header file named `portals.h`, although implementations are free to use a symlink to provide one or both of the files. Implementations should also provide a README file that explains implementation specific details. For example, it should list the limits (Section 3.6.1) for this implementation and provide a list of status registers that are provided (Section 3.3.7). See Appendix B for a template.

Numerous data structures are described as C-style structures in the Portals API; however, the definition is not meant to specify a field ordering. The implementation is free to optimize the ordering of data structures.

3.2. Constants

The Portals API defines a number of constants. Constants defined in this specification must be compile time constants. Further, constants whose type is specified to be integral must be valid labels for switch statements. Constants are generally associated with a base type in which constants are stored. Implementations are given freedom regarding the numeric values used for constants and their associated base types, constrained only by the compile time requirements.

3.2.1. Version Information

Every Portals implementation must provide two preprocessor constants, `PTL_MAJOR_VERSION` and `PTL_MINOR_VERSION`, which indicate the version of the Portals specification implemented by the implementation. In the case of versions that have two minor numbers the `PTL_MINOR_VERSION` will be equal to the first integer value of the minor version numbering (e.g. Portals 4.0.2 `PTL_MINOR_VERSION` will return 0).

Discussion: *`PTL_MAJOR_VERSION` and `PTL_MINOR_VERSION` were added to the Portals 4.0.2 specification. Previous versions, including Portals 4.0.1, did not include version constants. Users of the Portals interface must take that into account when using the version constants.*

3.3. Base Types

The Portals API defines a variety of base types. These types represent a simple renaming of the base types provided by the C programming language. In most cases these new type names have been introduced to improve type safety and to avoid issues arising from differences in representation sizes (e.g., 16-bit or 32-bit integers). Table 3-5 on page 114 lists all the types defined by Portals.

3.3.1. Sizes

The type `ptl_size_t` is an unsigned 64-bit integral type used for representing sizes. The constant `PTL_SIZE_MAX` represents the largest value a `ptl_size_t` can hold.

3.3.2. Handles

Objects maintained by the API are accessed through handles. Handle types have names of the form `ptl_handle_xx_t`, where `xx` is one of the two letter object type codes shown in Table 3-1, column `xx`. For example, the type `ptl_handle_ni_t` is used for network interface handles. Like all Portals types, their names use lower case letters and underscores are used to separate words.

Each type of object is given a unique handle type to enhance type checking. The type `ptl_handle_any_t` can be used when a generic handle is needed. Every handle value can be converted into a value of type `ptl_handle_any_t` without loss of information.

The type of a handle is left unspecified, but must be assignable in C. Every Portals object is associated with a specific network interface and the network handle associated with an object's handle may be retrieved by calling `PtINiHandle()`.

**IMPLEMENTATION
NOTE 4:**

Size of handle types

It is highly recommended that a handle type should be no larger than the native machine word size.

The constant `PTL_EQ_NONE`, of type `ptl_handle_eq_t`, is used to indicate the absence of an event queue. Similarly, the constant `PTL_CT_NONE`, of type `ptl_handle_ct_t`, indicates the absence of a counting event. See Section 3.10.1 for uses of these values. The special constant `PTL_INVALID_HANDLE` is used to represent an invalid handle.

**IMPLEMENTATION
NOTE 5:**

Unique handles

The encoding of handles is not specified by the Portals API. An implementation may reuse handle values, however the implementation is responsible for handling race conditions between threads calling release and acquire functions (such as **PtIMDRelease()** and **PtIMDBind()**).

3.3.3. Indexes

The type **ptl_pt_index_t** is an integral type used for representing portal table indexes. See Section 3.6.1 and 3.6.2 for limits on values of this type.

3.3.4. Match Bits

The type **ptl_match_bits_t** is capable of holding unsigned 64-bit integer values.

3.3.5. Network Interfaces

The type **ptl_interface_t** is an integral type used for identifying different network interfaces. Users will need to consult the implementation's README documentation to determine appropriate values for the interfaces available. The special constant **PTL_IFACE_DEFAULT** identifies the default interface.

3.3.6. Identifiers

The type **ptl_nid_t** is an integral type used for representing node identifiers and **ptl_pid_t** is an integral type for representing process identifiers when physical addressing is used in the network interface (**PTL_NI_PHYSICAL** is set for the network interface). If **PTL_NI_LOGICAL** is set, a *rank* (**ptl_rank_t**) is used instead. **ptl_uid_t** is an integral type for representing usage identifiers.

The special values **PTL_PID_ANY** matches any process identifier, **PTL_NID_ANY** matches any node identifier, **PTL_RANK_ANY** matches any rank, and **PTL_UID_ANY** matches any usage identifier. See Section 3.11 and 3.12 for uses of these values.

3.3.7. Status Registers

Each network interface maintains an array of status registers that can be accessed using the **PtINIStatus()** function (Section 3.6.4). The type **ptl_sr_index_t** defines the type of indexes that can be used to access the status registers. A small number of indexes are defined for all implementations:

Status Register Indexes (`ptl_sr_index_t`)

<code>PTL_SR_DROP_COUNT</code>	Identifies the status register that counts the dropped requests for the interface.
<code>PTL_SR_PERMISSION_VIOLATIONS</code>	Counts the number of attempted permission violations.
<code>PTL_SR_OPERATION_VIOLATIONS</code>	Counts the number of attempted operation violations

A permission violation is a violation of the usage id check, while an operation violation is a violation of the allowed operation types (put and/or get). Note that these three operations are orthogonal such that permission violations and operations violations should not increment `PTL_SR_DROP_COUNT`. Other indexes (and registers) may be defined by the implementation.

The type `ptl_sr_value_t` defines the type of values held in status registers. This is a signed integer type. The size is implementation dependent but must be at least 32 bits.

3.4. Function Arguments and Return Codes

Unless otherwise noted, an implementation is not required to check the validity of any arguments to a Portals function call. The argument to many Portals functions is a pointer to a type (because the argument is a pointer to a structure and/or because the argument is an output parameter). Unless otherwise noted, a pointer must point to a valid instance of the specified type; NULL is not generally a valid argument.

The Portals API specifies return codes that indicate success or failure of a function call. In the case where the failure is due to invalid arguments being passed into the function, the exact behavior of an implementation is undefined. The API suggests error codes that provide more detail about specific invalid parameters, but an implementation is not required to return these specific error codes. For example, an implementation is free to allow the caller to fault when given an invalid address, rather than return `PTL_ARG_INVALID`. In addition, an implementation is free to map these return codes to standard return codes where appropriate. For example, a Linux kernel-space implementation could map portals return codes to POSIX-compliant return codes. Table 3-7 on page 118 lists all return codes used by Portals.

3.5. Initialization and Cleanup

The Portals API includes a function, `PtlInit()`, to initialize the library and a function, `PtlFini()`, to clean up after the process is done using the library. The initialization state of Portals is reference counted so that repeated calls to `PtlInit()` and `PtlFini()` within a process (collection of threads) do not invalidate Portals state until the reference count reaches zero. Portals is *initialized* upon successful completion of the first call to `PtlInit()` and *finalized* upon successful completion of the first call to `PtlFini()` that results in the reference count reaching zero.

The Portals API also provides the `PtlAbort()` function, which allows an application to properly abort waiting and polling Portals routines.

A child process does not inherit any Portals resources from its parent. A child process must initialize Portals in order to obtain new, valid Portals resources. If a child process fails to initialize Portals and then uses the Portals interface, behavior is undefined for both the parent and the child.

3.5.1. PtlInit

The **PtlInit()** function initializes the Portals library. **PtlInit()** must be called at least once by a process before any thread makes a Portals function call and may be safely called more than once. Each call to **PtlInit()** increments a reference count. **PtlInit()** cannot be called after the Portals library has been finalized.

Function Prototype for PtlInit

```
int PtlInit(void);
```

Return Codes

PTL_OK	Indicates success.
PTL_FAIL	Indicates an error during initialization.

3.5.2. PtlFini

The **PtlFini()** function allows an application to clean up after the Portals library is no longer needed by a process. Each call to **PtlFini()** decrements the reference count that was incremented by **PtlInit()**. When the reference count reaches zero, all Portals resources are freed. Once the Portals resources are freed, calls to any of the functions defined by the Portals API or use of the structures set up by the Portals API will result in undefined behavior. Each call to **PtlInit()** should be matched by a corresponding **PtlFini()**.

Function Prototype for PtlFini

```
void PtlFini(void);
```

3.5.3. PtlAbort

The **PtlAbort()** function allows an application to gracefully abort the execution of waiting and polling Portals routines. Multithreaded applications with threads waiting in **PtlEQPoll()**, **PtlEQWait()**, **PtlCTWait()** or **PtlCTPoll()** must use **PtlAbort()** before calling **PtlEQFree()** and **PtlCTFree()** in order to avoid use-after-free scenarios.

Polling and waiting functions must stop and return once **PtlAbort()** is called. Subsequent calls to polling and waiting routines must immediately return **PTL_ABORTED**. Once **PtlAbort()** is completed, the application may safely call **PtlEQFree()** and **PtlCTFree()**.

Function Prototype for PtlAbort

```
void PtlAbort(void);
```

Discussion: *The use of **PtlAbort()** is not limited to multithreaded applications. Single-threaded programs may use **PtlAbort()** in signal handlers to properly terminate waiting and polling Portals functions. Notice also that **PtlAbort()** does not affect the behavior of **PtlCTGet()** and **PtlEQGet()** as they are non-blocking functions.*

3.6. Network Interfaces

The Portals API supports the use of multiple network interfaces. However, each interface is treated as an independent entity. Combining interfaces (e.g., “bonding” to create a higher bandwidth connection) must be handled internally by the Portals implementation, embedded in the underlying network, or handled by the application. Interfaces are treated as independent entities to make it easier to cache information on individual network interface cards.

A Portals *physical* network interface is a per-process abstraction of a physical network interface (or group of hardware interfaces). A physical network interface can not be used directly, but can be used by a process to instantiate up to four *logical* network interfaces. All logical network interfaces associated with a single physical network interface share the same network id and process id (nid/pid), but all other resources are unique to a logical network interface. A logical network interface can be initialized to provide either matching or non-matching Portals addressing and either logical or physical addressing of network endpoints through the data movement calls. These two options are independent and all four combinations of logical network interface options must be supported by each physical network interface.

Once initialized, each logical interface provides a portal table and a collection of status registers. In order to facilitate the development of portable Portals applications, a compliant implementation must provide at least 250 portal table entries. See Section 3.6.4 for a discussion of the **PtINIStatus()** function, which can be used to read the value of a status register. Every other type of Portals object (e.g., memory descriptor, event queue, or list entry) is also associated with a specific logical network interface. The association to a logical network interface is established when the object is created, and the **PtINIHandle()** function (Section 3.6.5) may be used to determine the logical network interface with which an object is associated.

Each logical network interface is initialized and shut down independently. The initialization routine, **PtINIInit()**, returns an interface object handle which is used in all subsequent portals operations. The **PtINIFini()** function is used to shut down a logical interface and release any resources that are associated with the interface. Network interface handles are associated with processes, not threads. All threads in a process share all of the network interface handles.

3.6.1. The Network Interface Limits Type

The function **PtINIInit()** accepts a pointer to a structure of desired limits and can fill a structure with the actual values supported by the network interface. Resource limits are specified independently for each logical network interface and the resources are shared by all users of the same network interface. The two structures are of type **ptl_ni_limits_t** and include the following members:


```

typedef struct {
    int max_entries;
    int max_unexpected_headers;
    int max_mds;
    int max_cts;
    int max_eqs;
    int max_pt_index;
    int max_iovecs;
    int max_list_size;
    int max_triggered_ops;
    ptl_size_t max_msg_size;
    ptl_size_t max_atomic_size;
    ptl_size_t max_fetch_atomic_size;
    ptl_size_t max_waw_ordered_size;
    ptl_size_t max_war_ordered_size;
    ptl_size_t max_volatile_size;
    unsigned int features;
} ptl_ni_limits_t;

```

Limits

<i>max_entries</i>	Maximum number of match list entries or list entries that can be allocated at any one time (only one of the two exists on an interface).
<i>max_unexpected_headers</i>	Maximum number of unexpected headers that the implementation can buffer.
<i>max_mds</i>	Maximum number of memory descriptors that can be allocated at any one time.
<i>max_eqs</i>	Maximum number of event queues that can be allocated at any one time.
<i>max_cts</i>	Maximum number of counting events that can be allocated at any one time.
<i>max_pt_index</i>	Largest portal table index for this interface, valid indexes range from 0 to <i>max_pt_index</i> , inclusive. An interface must support a <i>max_pt_index</i> of at least 249.
<i>max_iovecs</i>	Maximum number of I/O vectors for a single memory descriptor, list entry, or match list entry for this interface.
<i>max_list_size</i>	Maximum number of entries that can be attached to the list on any portal table index.
<i>max_triggered_ops</i>	Maximum number of triggered operations that can be outstanding.
<i>max_msg_size</i>	Maximum size (in bytes) of a message (<i>put</i> , <i>get</i> , or <i>reply</i>).
<i>max_atomic_size</i>	Maximum size (in bytes) that can be passed to an atomic operation. Any byte within an operation that is less than <i>max_atomic_size</i> is guaranteed to only be written to the user memory buffer once.
<i>max_fetch_atomic_size</i>	Maximum size (in bytes) that can be passed to an atomic operation that returns the prior value to the initiator.
<i>max_waw_ordered_size</i>	Maximum size (in bytes) of a message that will guarantee “per-address” data ordering for a write followed by a write (consecutive <i>put</i> or <i>atomic</i> or a mixture of the two) and a write followed by a read (<i>put</i> followed by a <i>get</i>). An interface must provide a <i>max_waw_ordered_size</i> of at least 64 bytes.

<i>max_war_ordered_size</i>	Maximum size (in bytes) of a message that will guarantee “per-address” data ordering for a read followed by a write (<i>get</i> followed by a <i>put</i> or <i>atomic</i>). An interface must provide a <i>max_war_ordered_size</i> of at least 8 bytes.
<i>max_volatile_size</i>	Maximum size (in bytes) that can be passed as the <i>length</i> of a <i>put</i> or <i>atomic</i> for a memory descriptor with the PTL_MD_VOLATILE option set.
<i>features</i>	A bit mask of features supported by the the Portals implementation. Currently, three features are defined. PTL_TARGET_BIND_INACCESSIBLE is discussed in Section 3.11 and 3.12, PTL_TOTAL_DATA_ORDERING is discussed in Section 2.7, and PTL_COHERENT_ATOMICS is discussed in Section 3.15.4.

3.6.2. PtlNIInit

The **PtlNIInit()** function initializes the Portals API for a network interface (NI). A process using Portals must call this function at least once before any other functions that apply to that interface. An additional call to **PtlSetMap()** must be made before communication calls are made on a logically addressed interface (See Section 3.6.6). Calls to **PtlNIInit()** increment a reference count on the network interface and must be matched by a call to **PtlNIInit()**. If **PtlNIInit()** gets called more than once *per logical interface*, then the implementation should fill in *actual* and *ni_handle* with the values obtained by the first caller and should ignore the *pid* argument. **PtlGetId()** or **PtlGetPhysId()** (Section 3.9) can be used to retrieve the *pid*.

Discussion: *Proper initialization of a logical network interface that uses logical endpoint addressing requires the user to call **PtlSetMap()**, creating a mapping of logical ranks to physical node IDs and process IDs. The physical address (NID/PID) associated with a logical network interface may be obtained by calling **PtlGetPhysId()**. The physical address may then be shared through an outside mechanism (including another Portals logical interface) to establish a consistent mapping of rank to NID/PID.*

Function Prototype for PtlNIInit

```
int PtlNIInit(ptl_interface_t iface,
             unsigned int options,
             ptl_pid_t pid,
             const ptl_ni_limits_t *desired,
             ptl_ni_limits_t *actual,
             ptl_handle_ni_t *ni_handle);
```

Arguments

<i>iface</i>	input	Identifies the physical network interface to be initialized. (See Section 3.3.5 for a discussion of values used to identify network interfaces.)
<i>options</i>	input	This field contains options that are requested for the network interface. Values for this argument can be constructed using a bitwise OR of the values defined below. Either PTL_NI_MATCHING or PTL_NI_NO_MATCHING must be set, but not both. Either PTL_NI_LOGICAL or PTL_NI_PHYSICAL must be set, but not both, to specify the endpoint addressing mode.

<i>pid</i>	input	Identifies the desired process identifier (for well known process identifiers). The specified <i>pid</i> must either be non-negative and less than the value <code>PTL_PID_MAX</code> or be <code>PTL_PID_ANY</code> . The value <code>PTL_PID_ANY</code> may be used to let the Portals library select a process identifier. See Section 3.9 for more information on process identifiers.
<i>desired</i>	input	If not NULL, points to a structure that holds the desired limits. If NULL, either previously set limits or implementation defined defaults will be used.
<i>actual</i>	output	If not NULL, on successful return, the location pointed to by <i>actual</i> will hold the actual limits.
<i>ni_handle</i>	output	On successful return, this location will hold the interface handle.

options

<code>PTL_NI_MATCHING</code>	Request that the interface specified in <i>iface</i> be opened with matching enabled.
<code>PTL_NI_NO_MATCHING</code>	Request that the interface specified in <i>iface</i> be opened with matching disabled. <code>PTL_NI_MATCHING</code> and <code>PTL_NI_NO_MATCHING</code> are mutually exclusive.
<code>PTL_NI_LOGICAL</code>	Request that the interface specified in <i>iface</i> be opened with logical endpoint addressing (e.g., GASNet node and rank or SHMEM PE).
<code>PTL_NI_PHYSICAL</code>	Request that the interface specified in <i>iface</i> be opened with physical endpoint addressing (e.g., NID/PID). <code>PTL_NI_LOGICAL</code> and <code>PTL_NI_PHYSICAL</code> are mutually exclusive.

Return Codes

<code>PTL_OK</code>	Indicates success.
<code>PTL_NO_INIT</code>	Indicates that the Portals API has not been successfully initialized.
<code>PTL_ARG_INVALID</code>	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
<code>PTL_PID_IN_USE</code>	Indicates that <i>pid</i> is currently in use.
<code>PTL_NO_SPACE</code>	Indicates that <code>PtlNlInit()</code> was not able to allocate the memory required to initialize the interface.

Discussion: *Each interface has its own sets of limits. In implementations that support multiple interfaces, the limits passed to and returned by `PtlNlInit()` apply only to the interface specified in *iface*. However, the use of *desired* is implementation dependent and an implementation may choose to ignore the request or provide limits based on a previous request.*

The desired limits are used to offer a hint to an implementation as to the amount of resources needed, and the implementation returns the actual limits available for use. In the case where an implementation does not have any pre-defined limits, it is free to return the largest possible value permitted by the corresponding type (e.g., `INT_MAX`). A quality implementation will enforce the limits that are returned and take the appropriate action when limits are exceeded, such as using the `PTL_NO_SPACE` return code. The caller is permitted to use maximum values for the desired fields to indicate that the limit should be determined by the implementation. An implementation must provide at least the resources specified by *actual*, unless the implementation returned the largest possible value permitted by the corresponding type in which case the implementation may be restricted by another resource such as available

application memory or machine capabilities that are beyond the Portals implementation's control. For example, a user may request a value for *max_unexpected_headers*, which a Portals implementation may return a value in *actual* of the maximum value for that type (e.g. 64 bit integer). This indicates that the Portals implementation is able to dynamically allocate memory to support buffering unexpected headers. The Portals implementation will provide as much buffering as is possible for unexpected headers until it runs out of memory.

3.6.3. PtlNIFini

The **PtlNIFini()** function is used to release the resources allocated for a network interface. The release of network interface resources is based on a reference count that is incremented by **PtlNIInit()** and decremented by **PtlNIFini()**. Resources can only be released when the reference count reaches zero. Once the release of resources has begun, the results of pending API operations (e.g., operations initiated by another thread) for this interface are undefined. Similarly, the effects of incoming operations (*put*, *get*, *atomic*) or return values (*acknowledgment* and *reply*) for this interface are undefined until the interface is reinitialized by another call to **PtlNIInit()**.

Function Prototype for PtlNIFini

```
int PtlNIFini(ptl_handle_ni_t ni_handle);
```

Arguments

ni_handle **input** An interface handle to shut down.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.6.4. PtlNIStatus

The **PtlNIStatus()** function returns the value of a status register for the specified interface. See Section 3.3.7 for more information on status register indexes and status register values.

Function Prototype for PtlNIStatus

```
int PtlNIStatus(ptl_handle_ni_t ni_handle,  
                ptl_sr_index_t status_register,  
                ptl_sr_value_t *status);
```

Arguments

<i>ni_handle</i>	input	An interface handle.
<i>status_register</i>	input	The index of the status register.
<i>status</i>	output	On successful return, this location will hold the current value of the status register.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.6.5. PtlNIHandle

The **PtlNIHandle()** function returns the network interface handle with which the object identified by *handle* is associated. If the object identified by *handle* is a network interface, this function returns the same value it is passed.

Function Prototype for PtlNIHandle

```
int PtlNIHandle(ptl_handle_any_t handle,
               ptl_handle_ni_t *ni_handle);
```

Arguments

<i>handle</i>	input	The object handle.
<i>ni_handle</i>	output	On successful return, this location will hold the network interface handle associated with <i>handle</i> .

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.6.6. PtlSetMap

The **PtlSetMap()** function initializes the mapping from *logical* endpoint identifiers (rank) to *physical* endpoint identifiers (nid/pid) for the given logically addressed logical network interface. A process must ensure that the logical mapping is set before the specified logically addressed logical network interface may be used in any portals calls other than **PtlNIInit()**, **PtlGetMap()**, and **PtlGetPhysId()**. If the map of the other logically addressed logical network interface associated with the same physical network interface as the specified interface handle has not been set by a

call to **PtlSetMap()**, the implementation may choose to set the mapping on both logical network interfaces. It is erroneous to call **PtlSetMap()** on a physically addressed logical network interface. Subsequent calls (either by different threads or the same thread) to **PtlSetMap()** will overwrite any mapping associated with the logical network interface; hence, libraries must take care to ensure reasonable interoperability.

Function Prototype for PtlSetMap

```
int PtlSetMap(ptl_handle_ni_t ni_handle,
             ptl_size_t map_size,
             const ptl_process_t *mapping);
```

Arguments

<i>ni_handle</i>	input	The interface handle identifying the network interface which should be initialized with <i>mapping</i> . The network interface handle must refer to a logically addressed network interface.
<i>map_size</i>	input	The number of elements in <i>mapping</i> .
<i>mapping</i>	input	Points to an array of ptl_process_t structures where entry N in the array contains the NID/PID pair that is associated with the logical rank N.

Return Codes

PTL_OK	Indicates success.
PTL_IGNORED	Indicates no error occurred, but the implementation does not support dynamic changing of the logical identifier map, likely due to integration with a static run-time system.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_SPACE	Indicates that PtlSetMap() was not able to allocate the memory required to initialize the map.

Discussion: *PtlSetMap()* is a local operation and the map set by different communicating processes may be different. The *rank* field of target-side events may be unexpected in cases where the two processes have different maps.

3.6.7. PtlGetMap

The **PtlGetMap()** function retrieves the mapping from *logical* identifiers (rank) to *physical* identifiers (nid/pid) for the specified logically addressed logical network interface. If the *map_size* is smaller than the actual map size, the first *map_size* entries in the map will be copied into *mapping*. If the *map_size* is larger than the actual map size, the entire map is copied into *mapping* and the buffer beyond the *actual_map_size* entry is left unmodified. It is erroneous to call **PtlGetMap()** on a physically addressed logical network interface.

Function Prototype for PtlGetMap

```
int PtlGetMap(ptl_handle_ni_t ni_handle,
              ptl_size_t map_size,
              ptl_process_t *mapping,
              ptl_size_t *actual_map_size);
```

Arguments

<i>ni_handle</i>	input	The network interface handle from which the map should be retrieved. The network interface handle must refer to a logically addressed logical network interface.
<i>map_size</i>	input	The length of <i>mapping</i> in number of elements.
<i>mapping</i>	output	Points to an array of ptl_process_t structures where entry N in the array will be populated with the NID/PID pair that is associated with the logical rank N.
<i>actual_map_size</i>	output	On return, <i>actual_map_size</i> contains the size, in number of elements, of the map currently associated with the logical interface. May be bigger than <i>map_size</i> or the <i>mapping</i> array.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_IGNORED	Indicates that the request was ignored as there was no map set on the logical network interface.

3.7. Portal Table Entries

A portal index refers to a portal table entry. The assignment of these indexes can either be statically or dynamically managed, and will typically be a combination of both. A portal table entry must be allocated before being used. From a user perspective, messages that arrive traverse list entries or match list entries in the order they were appended within a single portal table index. Resource exhaustion (Section 2.8) is handled independently on different portal table entries.

3.7.1. PtlPTAlloc

The **PtlPTAlloc()** function allocates a portal table entry and sets flags that pass options to the implementation. A portal table entry allocated with **PTL_PT_ALLOC_DISABLED** is initialized to disabled status (see Section 3.7.3), and no new messages should be accepted on the specified portal table entry until enabled with **PtlPTEnable()** (Section 3.7.4).

Function Prototype for PtlPTAlloc

```
int PtlPTAlloc(ptl_handle_ni_t ni_handle,
               unsigned int options,
               ptl_handle_eq_t eq_handle,
               ptl_pt_index_t pt_index_req,
               ptl_pt_index_t *pt_index);
```

Arguments

<i>ni_handle</i>	input	The interface handle to use.
<i>options</i>	input	This field contains options that are requested for the portal index. Values for this argument can be constructed using a bitwise OR of the values defined below.
<i>eq_handle</i>	input	The event queue handle used to log the events related to the list entries attached to the portal table entry. If this argument is PTL_EQ_NONE, events related to this portal table entry are not logged.
<i>pt_index_req</i>	input	The value of the portal index that is requested. If the value is set to PTL_PT_ANY, the implementation can return any portal index.
<i>pt_index</i>	output	On successful return, this location will hold the portal index that has been allocated.

Options

PTL_PT_ONLY_USE_ONCE	Indicate to the underlying implementation that all entries attached to the priority list on this portal table entry are guaranteed to have the PTL_ME_USE_ONCE or PTL_LE_USE_ONCE option set.
PTL_PT_ONLY_TRUNCATE	Indicate to the underlying implementation that all entries attached to the priority list on this portal table entry are guaranteed not to have the PTL_ME_NO_TRUNCATE option set.
PTL_PT_FLOWCTRL	Enable flow control on this portal table entry (see Section 2.8).
PTL_PT_ALLOC_DISABLED	Allocate this portal table entry in disabled status (see Section 3.7.3).

Return Codes

PTL_OK	Indicates success.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_PT_FULL	Indicates that there are no free entries in the portal table.
PTL_PT_IN_USE	Indicates that the Portal table entry requested is in use.
PTL_PT_EQ_NEEDED	Indicates that flow control is enabled and there is no EQ attached.

Discussion: The `PTL_PT_ONLY_USE_ONCE` and `PTL_PT_ONLY_TRUNCATE` options are hints to the implementation that convey that the user will be employing certain common usage scenarios when using the priority list. Use of these options may allow the implementation to optimize the matching logic. Note that the optimal set of options may vary depending on whether matching or non-matching logical network interfaces are used. For a matching logical network interface, an implementation likely may optimize the case where both `PTL_PT_ONLY_USE_ONCE` and `PTL_PT_ONLY_TRUNCATE` are specified. For a non-matching logical network interface, pre-posted persistent LEs are likely to provide better performance.

Discussion: To provide an opportunity to guarantee resources are available (and hence avoid triggering flow control on entries for which flow control is enabled), it is recommended that users allocate portal table entries with the `PTL_PT_ALLOC_DISABLED` option, and subsequently explicitly enable the entry using `PtIPTEnable()`.

3.7.2. PtIPTFree

The `PtIPTFree()` function releases the resources associated with a portal table entry. Objects associated with the portal table entry, such as list entries and event queues, are not freed as the result of a call to `PtIPTFree()`.

Function Prototype for `PtIPTFree`

```
int PtIPTFree( ptl_handle_ni_t ni_handle,
              ptl_pt_index_t pt_index);
```

Arguments

- | | | |
|------------------|--------------|--|
| <i>ni_handle</i> | input | The interface handle on which the <i>pt_index</i> should be freed. |
| <i>pt_index</i> | input | The portal index that is to be freed. |

Return Codes

- | | |
|------------------------------|---|
| <code>PTL_OK</code> | Indicates success. |
| <code>PTL_NO_INIT</code> | Indicates that the Portals API has not been successfully initialized. |
| <code>PTL_ARG_INVALID</code> | Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent. |
| <code>PTL_PT_IN_USE</code> | Indicates that <i>pt_index</i> is currently in use (e.g., a match list entry is still attached). |

3.7.3. PtIPTDisable

The `PtIPTDisable()` function indicates to an implementation that no new messages should be accepted on the specified portal table entry. The function blocks until the portal table entry status has been updated, all messages being actively processed are completed, and all events are delivered. Since `PtIPTDisable()` waits until the portal table entry is disabled before it returns, it does not generate a `PTL_EVENT_PT_DISABLED` event. Processing of operations targeting other portal table entries and local operations continues after a call to `PtIPTDisable()`.

Function Prototype for PtlPTDisable

```
int PtlPTDisable(ptl_handle_t ni_handle,  
                ptl_pt_index_t pt_index);
```

Arguments

<i>ni_handle</i>	input	The interface handle to use.
<i>pt_index</i>	input	The portal index that is to be disabled.

Return Codes

PTL_OK	Indicates success.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.

3.7.4. PtlPTEnable

The **PtlPTEnable()** function indicates to an implementation that a previously disabled portal table entry should be re-enabled. This is used to enable portal table entries that were automatically or manually disabled. The function blocks until the portal table entry is enabled.

Function Prototype for PtlPTEnable

```
int PtlPTEnable(ptl_handle_t ni_handle,  
               ptl_pt_index_t pt_index);
```

Arguments

<i>ni_handle</i>	input	The interface handle to use.
<i>pt_index</i>	input	The value of the portal index to enable.

Return Codes

PTL_OK	Indicates success.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.

Discussion: *PtlPTEnable()* re-enables a portal table entry, allowing incoming messages to match against list entries associated with the portal table entry. Messages may have been dropped while the portal table entry was disabled. Higher level communication protocols with strict ordering constraints may have to quiesce messages and retransmit after re-enabling a portal table entry (See Section 2.8).

3.8. Usage Identification

A usage identifier (UID) is assigned to a process through an implementation-specific mechanism. UIDs define access control with respect to remote Portals network interfaces. The usage identifier is included in a trusted header of a portals message, such that the trusted header is not modifiable by the user. They can be used at the *target* to limit access to list entries (Section 3.11 and Section 3.12). The UID is common across logical network interfaces within the same process, even if the logical network interfaces are over different physical network interfaces.

3.8.1. PtlGetUid

The `PtlGetUid()` function is used to retrieve the usage identifier of a process.

Function Prototype for PtlGetUid

```
int PtlGetUid(ptl_handle_ni_t ni_handle,  
             ptl_uid_t *uid);
```

Arguments

<i>ni_handle</i>	input	A network interface handle.
<i>uid</i>	output	On successful return, this location will hold the usage identifier for the calling process.

Return Codes

<code>PTL_OK</code>	Indicates success.
<code>PTL_ARG_INVALID</code>	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
<code>PTL_NO_INIT</code>	Indicates that the Portals API has not been successfully initialized.

3.9. Process Identification

Processes that use the Portals API can be identified using a node identifier and process identifier. Every node accessible through a network interface has a unique node identifier and every process running on a node has a unique process identifier. As such, any process in the computing system can be uniquely identified by its node identifier and process identifier. The node identifier and process identifier can be aggregated by the application into a rank, which is translated by the implementation into a network identifier and process identifier. It is an implementation decision whether two physical network interfaces in the same node have the same node or process identifiers. All logical network interfaces which share the same physical network interface share the same node and process identifiers.

The Portals API defines a type, `ptl_process_t`, for representing process identifiers, and two functions, `PtlGetId()` and `PtlGetPhysId()`, which can be used to obtain the identifier of the current process.

Discussion: *The Portals API does not include thread identifiers. Messages are delivered to processes (address spaces) not threads (contexts of execution).*

3.9.1. The Process Identification Type

The `ptl_process_t` type is a union that can represent the process as either a physical address or a logical address within the machine. The physical address uses two identifiers to represent a process identifier: a node identifier *nid* and a process identifier *pid*. In turn, a logical address uses a logical index within a translation table specified by the application (the *rank*) to identify another process.

```
typedef union {
    struct {
        ptl_nid_t nid;
        ptl_pid_t pid;
    } phys;
    ptl_rank_t rank;
} ptl_process_t;
```

3.9.2. PtlGetId

Function Prototype for `PtlGetId`

```
int PtlGetId(ptl_handle_ni_t ni_handle,
             ptl_process_t *id);
```

Arguments

<i>ni_handle</i>	input	A network interface handle.
<i>id</i>	output	On successful return, this location will hold the identifier for the calling process. If the interface is logically addressed, the logical address is returned. If the interface is physically addressed, the physical address is returned.

Return Codes

PTL_OK	Indicates success.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.

Discussion: *Note that process identifiers and ranks are dependent on the network interface(s). In particular, if a node has multiple interfaces, it may have multiple process identifiers and multiple ranks.*

3.9.3. PtlGetPhysId

Function Prototype for PtlGetPhysId

```
int PtlGetPhysId(ptl_handle_ni_t ni_handle,  
                ptl_process_t *id);
```

Arguments

<i>ni_handle</i>	input	A network interface handle.
<i>id</i>	output	On successful return, this location will hold the identifier for the calling process. The physical address is always returned, even for logically addressed network interfaces.

Return Codes

PTL_OK	Indicates success.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.

Discussion: Note that process identifiers and ranks are dependent on the network interface(s). In particular, if a node has multiple interfaces, it may have multiple process identifiers and multiple ranks.

3.10. Memory Descriptors

A memory descriptor contains information about a region of a process' memory and optionally points to an event queue and counting event where information about the operations performed on the memory descriptor are recorded. Memory descriptors are initiator side resources that are used to encapsulate the association of a network interface (NI) with a description of a memory region. They provide an interface to register memory (for operating systems that require it) and to carry that information across multiple operations (an MD is persistent until released). **PtlMDBind()** is used to create a memory descriptor and **PtlMDRelease()** is used to unlink and release the resources associated with a memory descriptor.

A memory descriptor describes a memory region using a base address and length; however, it is not a requirement for all of the memory described by the memory descriptor to be allocated or accessible within the application. For example, an application can create a memory descriptor that covers the entire virtual address range by setting *start* to NULL and *length* to PTL_SIZE_MAX, even though the entire region is not currently allocated. If the application issues a portals operation (e.g., *put*) that would access an unallocated region of the MD, the implementation may either cause a segmentation fault of the application or may simply fail the operation. If a full event is delivered, it must set *ni_fail_type* to PTL_NI_SEGV. If the memory descriptor sets the PTL_IOVEC option, the memory region(s) described by the **ptl_iovec_t** must all be accessible within the application.

**IMPLEMENTATION
NOTE 6:**

Memory descriptors that bind inaccessible memory

The implementation is responsible for handling any issues, such as the memory registration required by some platforms, that arise from the ability of an MD to cover all of the virtual address space. While some implementations may have elegant solutions to this issue (e.g., lightweight kernels or NIC hardware translation caching), other implementations may require registration caching schemes.

3.10.1. The Memory Descriptor Type

The `ptl_md_t` type defines the visible parts of a memory descriptor. Values of this type are used to initialize the memory descriptors.

```
typedef struct {  
    void *start;  
    ptl_size_t length;  
    unsigned int options;  
    ptl_handle_eq_t eq_handle;  
    ptl_handle_ct_t ct_handle;  
} ptl_md_t;
```

Members

start, length

Specify the memory region associated with the memory descriptor. The *start* member specifies the starting address for the memory region and the *length* member specifies the length of the region. There are no restrictions on buffer alignment, the starting address or the length of the region; although messages that are not natively aligned (e.g., to a four byte or eight byte boundary) may be slower (i.e., lower bandwidth and/or longer latency) on some implementations.

options

Specifies the behavior of the memory descriptor. Options include the use of scatter/gather vectors and control of events associated with this memory descriptor. Values for this argument can be constructed using a bitwise OR of the following values:

PTL_MD_EVENT_SEND_DISABLE

Specifies that this memory descriptor should not generate send events (PTL_EVENT_SEND). This flag does not affect counting events.

PTL_MD_EVENT_SUCCESS_DISABLE

Specifies that this memory descriptor should not generate full events if the *ni_fail_type* would be **PTL_OK**. This flag does not affect counting events. Disabling full events for successful operations is useful in scenarios when a counting event is sufficient for completion, but more information is needed for error recovery.

PTL_MD_EVENT_CT_SEND

Enable the counting of PTL_EVENT_SEND events.

PTL_MD_EVENT_CT_REPLY

Enable the counting of PTL_EVENT_REPLY events.

PTL_MD_EVENT_CT_ACK

Enable the counting of PTL_EVENT_ACK events.

PTL_MD_EVENT_CT_BYTES	By default, counting events count events. When set, this option causes bytes to be counted instead for success events. Byte counts must be incremented exactly once per operation. The increment is by the <i>mlength</i> that would be specified by the associated full event. Failure events always increment the count by one.
PTL_MD_UNORDERED	Indicate to the Portals implementation that messages sent from this memory descriptor do not have to arrive at the target in order. Note that this has no impact on acknowledgments or replies, which are never required to be ordered.
PTL_MD_VOLATILE	Indicate to the Portals implementation that the application may modify any send buffers associated with this memory descriptor immediately following the return from a portals operation if the <i>length</i> argument is less than or equal to <i>max_volatile_size</i> . In that case, operations should not return until it is safe for the application to reuse any send buffers. Note that the MD can be of any size, but the Portals implementation must honor this option as long as the operation (e.g., <i>put</i> or <i>atomic</i> , not <i>get</i>) uses a <i>length</i> less than or equal to <i>max_volatile_size</i> . Operations with <i>length</i> greater than <i>max_volatile_size</i> may not honor PTL_MD_VOLATILE, but should not return an error solely because the <i>length</i> is greater than <i>max_volatile_size</i> .
PTL_MD_UNRELIABLE	Indicate to the Portals implementation that guaranteed delivery of messages is not required and messages may be dropped by the network. Only the put operation with no acknowledgment is supported in this mode. Other operations are considered undefined behavior.
PTL_IOVEC	Specifies that the <i>start</i> argument is a pointer to an array of type ptl_iovec_t (Section 3.10.2) and the <i>length</i> argument is the length of the array of ptl_iovec_t elements. This allows for a scatter/gather capability for memory descriptors. A scatter/gather memory descriptor behaves exactly as a memory descriptor that describes a single virtually contiguous region of memory. The array of ptl_iovec_t elements referred to by the <i>start</i> argument cannot be changed or released for the lifetime of the memory descriptor.
<i>eq_handle</i>	The event queue handle used to log the operations performed on the memory region. If this argument is PTL_EQ_NONE, operations performed on this memory descriptor are not logged.
<i>ct_handle</i>	A handle for counting events associated with the memory region. If this argument is PTL_CT_NONE, operations performed on this memory descriptor are not counted.

3.10.2. The I/O Vector Type

The **ptl_iovec_t** type is used to describe scatter/gather buffers of a memory descriptor, list entry, or match list entry in conjunction with the PTL_IOVEC option. The **ptl_iovec_t** type is intended to be a type definition of the `struct iovec` type on systems that already support this type.

The **ptl_iovec_t** array is passed as the *start* field when creating a memory descriptor, list entry, or match list entry. It must not be modified or destroyed by the application or implementation for the life of the descriptor or entry. Descriptors or entries using **ptl_iovec_t** types may be combined with offsets (local and remote). The offset is computed as if the region described by the **ptl_iovec_t** type were a single contiguous region.

Discussion: Performance conscious users should not mix offsets (local or remote) with *ptl_iovec_t*. While this is a supported operation, it may have unexpected performance consequences.

```
typedef struct {
    void *iov_base;
    ptl_size_t iov_len;
} ptl_iovec_t;
```

Members

iov_base The byte aligned start address of the vector element
iov_len The length (in bytes) of the vector element

3.10.3. PtlMDBind

The **PtlMDBind()** operation is used to create a memory descriptor to be used by the *initiator*.

Function Prototype for PtlMDBind

```
int PtlMDBind(ptl_handle_ni_t ni_handle,
              const ptl_md_t *md,
              ptl_handle_md_t *md_handle);
```

Arguments

ni_handle **input** The network interface handle with which the memory descriptor will be associated.
md **input** Provides initial values for the user-visible parts of a memory descriptor. Other than its use for initialization, there is no linkage between this structure and the memory descriptor maintained by the implementation.
md_handle **output** On successful return, this location will hold the newly created memory descriptor handle. The *md_handle* argument must be a valid address and cannot be NULL.

Return Codes

PTL_OK Indicates success.
PTL_NO_INIT Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID Indicates that an invalid argument was passed. Argument checking is implementation dependent, but this may indicate that an invalid *ni_handle* was used, an invalid event queue was associated with the *md*, or other contents in the *md* were illegal.
PTL_NO_SPACE Indicates that there is insufficient memory to allocate the memory descriptor.

**IMPLEMENTATION
NOTE 7:**

Memory registration

On systems that require memory registration, the **PtIMDBind()** operation should invoke the appropriate memory registration functions.

**IMPLEMENTATION
NOTE 8:**

Optimization for Duplicate Memory Descriptors

Because the *eq_handle* and *ct_handle* are bound to the memory descriptor on the initiator, there are usage models where it is necessary to create numerous memory descriptors that only differ in their *eq_handle* or *ct_handle* field. Implementations may desire to optimize for this usage model.

3.10.4. PtIMDRelease

The **PtIMDRelease()** function releases the internal resources associated with a memory descriptor. (This function does not free the memory region associated with the memory descriptor; i.e., the memory the user allocated for this memory descriptor.) Only memory descriptors with no pending operations may be unlinked. A memory descriptor is considered to have pending operations if an operation has been started and the corresponding `PTL_EVENT_SEND` or `PTL_EVENT_REPLY` operation has not been delivered. A memory descriptor may be released before a `PTL_EVENT_ACK` event is delivered, in which case the acknowledgment will be discarded.

Function Prototype for PtIMDRelease

```
int PtIMDRelease(ptl_handle_md_t md_handle);
```

Arguments

md_handle **input** The memory descriptor handle to be released.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.11. List Entries and Lists

A list is a chain of list entries. Examples of lists include the priority list and the overflow list. Each list entry (LE) describes a memory region and includes a set of options. It is the target side analogue of the memory descriptor (MD) for non-matching logical network interfaces. The **PtILEAppend()** function appends a single list entry to the specified list on the specified portal index and returns the list entry handle. List entries can be dynamically removed from a list using the **PtILEUnlink()** function.

Like a memory descriptor, a list entry describes a memory region using a base address and length. A zero-length list entry may be created by setting *start* to NULL and *length* to 0. Zero-length buffers (NULL LE) are useful to record events. Messages that are outside the bounds of the LE are truncated to zero bytes (e.g., zero-length buffers or an offset beyond the length of the LE). If the interface set the PTL_TARGET_BIND_INACCESSIBLE bit in the *features* field of the *actual* limits (See Section 3.6.1), then it is not a requirement for all of the memory described by the list entry to be allocated or accessible within the application. For example, an application could create a list entry that covers the entire virtual address range by setting *start* to NULL and *length* to PTL_SIZE_MAX, even though the entire region is not currently allocated. If an incoming operation (e.g., *put*) attempts to access an unallocated region of the LE, the implementation may either cause a segmentation fault of the application or may simply fail the operation. If a full event is delivered, it must set *ni_fail_type* to PTL_NI_SEGV. The target may, however, set the PTL_LE_IS_ACCESSIBLE option to indicate that the entire memory space described by the LE is accessible. If the list entry sets the PTL_IOVEC option, the memory region(s) described by the *ptl_iovec_t* array must all be accessible within the application.

**IMPLEMENTATION
NOTE 9:**

List entries that bind inaccessible memory

If the implementation returns PTL_TARGET_BIND_INACCESSIBLE, then the implementation is responsible for handling any issues, such as the memory registration required by some platforms, that arise from the ability of an LE to cover all of the virtual address space. While some implementations may have elegant solutions to this issue (e.g., lightweight kernels or NIC hardware translation caching), other implementations may require a software thread on the target to implement a remote registration caching scheme like Firehose [3].

List entries can be appended to either the priority list or the overflow list associated with a portal table entry; however, when attached to an overflow list, additional semantics are implied that require the implementation to track messages that arrive in list entries. Essentially, the memory region identified is provided to the implementation for use in managing unexpected messages. Buffers provided in the overflow list will post a full event (PTL_EVENT_AUTO_UNLINK) when the buffer space has been consumed, to notify the application that more buffer space may be needed. When the application is free to reuse the buffer (i.e. the implementation is done with it), another full event (PTL_EVENT_AUTO_FREE) will be posted. The PTL_EVENT_AUTO_FREE full event will be posted after *all* other events (including counting events) associated with the buffer have been delivered.

Discussion: *It is the responsibility of the application to ensure that the implementation has sufficient buffer space to manage unexpected messages (i.e. in the unexpected list). Failure to do so will cause messages to be dropped. The PTL_EVENT_ACK at the initiator will indicate the failure as described in Section 3.13.3. Note that overflow events can readily exhaust the event queue. Proper use of the API will generally require the application to post at least two (and typically several) buffers so that the application has time to notice the PTL_EVENT_AUTO_UNLINK and replace the buffer. In many usage scenarios, however, the application may choose to have only persistent list entries—list entries without the PTL_LE_USE_ONCE option set—in the priority list. Thus, overflow list entries will not be required.*

It is the responsibility of the implementation to determine when a buffer that is automatically unlinked from an overflow list can be reused. It must note that it is no longer holding state associated with the buffer and post a PTL_EVENT_AUTO_FREE full event after all other events, including counting events, associated with that buffer have been delivered.

List entries can be appended to a network interface with the PTL_NI_NO_MATCHING option set (a non-matching network interface). A matching network interface requires a match list entry.

3.11.1. The List Entry Type

The `ptl_le_t` type defines the visible parts of a list entry. Values of this type are used to initialize the list entries.

```
typedef struct {
    void *start;
    ptl_size_t length;
    ptl_handle_ct_t ct_handle;
    ptl_uid_t uid;
    unsigned int options;
} ptl_le_t;
```

Members

start, length

Specify the memory region associated with the list entry. The *start* member specifies the starting address for the memory region and the *length* member specifies the length of the region. There are no restrictions on buffer alignment, the starting address or the length of the region; although messages that are not natively aligned (e.g., to a four byte or eight byte boundary) may be slower (i.e., lower bandwidth and/or longer latency) on some implementations.

ct_handle

A handle for counting events associated with the memory region. If this argument is `PTL_CT_NONE`, operations performed on this list entry are not counted.

uid

Specifies the usage ID that may access this list entry. The usage ID may be set to a wildcard (`PTL_UID_ANY`). If the access control check fails, then the message is dropped without modifying Portals state. This is treated as a permissions failure and the status register indexed by `PTL_SR_PERMISSION_VIOLATIONS` is incremented. This failure is also indicated to the initiator. If a full event is delivered to the initiator, the *ni_fail_type* in the `PTL_EVENT_ACK` event must be set to `PTL_NI_PERM_VIOLATION`.

options

Specifies the behavior of the list entry. The following options can be selected: enable *put* operations (yes or no), enable *get* operations (yes or no), message truncation (yes or no), acknowledgment (yes or no), use scatter/gather vectors and control event delivery. Values for this argument can be constructed using a bitwise OR of the following values:

`PTL_LE_OP_PUT`

Specifies that the list entry will respond to *put* operations. By default, list entries reject *put* operations. If a *put* operation targets a list entry where `PTL_LE_OP_PUT` is not set, it is treated as an operations failure and `PTL_SR_OPERATION_VIOLATIONS` is incremented. If a full event is delivered to the initiator, the *ni_fail_type* in the `PTL_EVENT_ACK` event must be set to `PTL_NI_OP_VIOLATION`.

`PTL_LE_OP_GET`

Specifies that the list entry will respond to *get* operations. By default, list entries reject *get* operations. If a *get* operation targets a list entry where `PTL_LE_OP_GET` is not set, it is treated as an operations failure and `PTL_SR_OPERATION_VIOLATIONS` is incremented. If a full event is delivered to the initiator, the *ni_fail_type* in the `PTL_EVENT_ACK` event must be set to `PTL_NI_OP_VIOLATION`.

	<p>Note: It is not considered an error to have a list entry that does not respond to either <i>put</i> or <i>get</i> operations: Nor is it considered an error to have a list entry that responds to both <i>put</i> and <i>get</i> operations. In fact, a list entry must be configured to respond to both <i>put</i> and <i>get</i> operations to properly handle a PtlFetchAtomic() or PtlSwap() operation.</p>
PTL_LE_USE_ONCE	Specifies that the list entry will only be used once and then unlinked. If this option is not set, the list entry persists until it is explicitly unlinked.
PTL_LE_UNEXPECTED_HDR_DISABLE	Specifies that the header for a message delivered to this list entry should not be added to the unexpected list. This option only has meaning if the list entry is inserted into the overflow list. By creating a list entry which truncates messages to zero bytes, disables comm events, and sets this option, a user may create a list entry which consumes no target side resources. A list entry with this flag set does not generate PTL_EVENT_AUTO_FREE events.
PTL_IOVEC	Specifies that the <code>start</code> argument is a pointer to an array of type ptl_iovec_t (Section 3.10.2) and the <code>length</code> argument is the length of the array. This allows for a scatter/gather capability for list entries. A scatter/gather list entry behaves exactly as a list entry that describes a single virtually contiguous region of memory. All other semantics are identical. The array of ptl_iovec_t elements referred to by the <code>start</code> argument cannot be changed or released until the list entry is unlinked.
PTL_LE_IS_ACCESSIBLE	Indicate that this list entry only contains memory addresses that are accessible by the application.
PTL_LE_EVENT_LINK_DISABLE	Specifies that this list entry should not generate a PTL_EVENT_LINK full event indicating the list entry successfully linked.
PTL_LE_EVENT_COMM_DISABLE	Specifies that this list entry should not generate full events that indicate a communication operation. This includes PTL_EVENT_GET, PTL_EVENT_PUT, PTL_EVENT_ATOMIC, PTL_EVENT_FETCH_ATOMIC, and PTL_EVENT_SEARCH.
PTL_LE_EVENT_FLOWCTRL_DISABLE	Specifies that this list entry should not generate a PTL_EVENT_PT_DISABLED full event indicating a flow control failure when the current list entry generated the failure.
PTL_LE_EVENT_SUCCESS_DISABLE	Specifies that this list entry should not generate full events if the <i>ni_fail_type</i> would be PTL_OK for the following events: PTL_EVENT_PUT, PTL_EVENT_GET, PTL_EVENT_ATOMIC, PTL_EVENT_FETCH_ATOMIC, PTL_EVENT_SEARCH, PTL_EVENT_PUT_OVERFLOW, PTL_EVENT_GET_OVERFLOW, PTL_EVENT_ATOMIC_OVERFLOW, and PTL_EVENT_FETCH_ATOMIC_OVERFLOW. This flag does not affect counting events. Disabling full events for successful operations is useful in scenarios when a counting event is sufficient for completion, but more information is needed for error recovery.
PTL_LE_EVENT_OVER_DISABLE	Specifies that this list entry should not generate overflow list full events. This includes PTL_EVENT_PUT_OVERFLOW, PTL_EVENT_GET_OVERFLOW, PTL_EVENT_ATOMIC_OVERFLOW, and PTL_EVENT_FETCH_ATOMIC_OVERFLOW.
PTL_LE_EVENT_UNLINK_DISABLE	Specifies that this list entry should not generate auto-unlink (PTL_EVENT_AUTO_UNLINK) or free (PTL_EVENT_AUTO_FREE) full events.
PTL_LE_EVENT_CT_COMM	Enable the counting of communication full events (PTL_EVENT_PUT, PTL_EVENT_GET, PTL_EVENT_ATOMIC, PTL_EVENT_FETCH_ATOMIC, and PTL_EVENT_SEARCH).

PTL_LE_EVENT_CT_OVERFLOW	Enable the counting of overflow events (PTL_EVENT_PUT_OVERFLOW, PTL_EVENT_GET_OVERFLOW, PTL_EVENT_ATOMIC_OVERFLOW, PTL_EVENT_FETCH_ATOMIC_OVERFLOW).
PTL_LE_EVENT_CT_BYTES	By default, counting events count events. When set, this option causes bytes to be counted instead for success events. Byte counts must be incremented exactly once per operation. The increment is by the number of bytes counted (<i>mlength</i>). Failure events always increment the count by one.

Discussion: When the `PTL_LE_USE_ONCE` option is set, an event associated with a target side operation (e.g., a `PTL_EVENT_PUT` full event) also implies that the associated list entry has unlinked; hence, it is safe on these list entries to set the `PTL_LE_EVENT_UNLINK_DISABLE` option.

`PTL_LE_EVENT_FLOWCTRL_DISABLE` only disables flow control events which are the direct result of an incoming message matching the current list entry. This includes a message matching the list entry but the associated event queue is full or a message matching a list entry in the overflow list but the unexpected headers list is full. If flow control is enabled on the portal table entry and a message does not match in either the priority or overflow lists, a `PTL_EVENT_PT_DISABLED` event is always generated.

3.11.2. PtlLEAppend

The `PtlLEAppend()` function creates a single list entry and appends this entry to the end of the list specified by `ptl_list` associated with the portal table entry specified by `pt_index` for the portal table for `ni_handle`.

When a list entry is posted to a priority list, the unexpected list is checked to see if a message has arrived prior to posting the list entry. If so, an appropriate overflow full event is generated, the matching header is removed from the unexpected list, and a list entry with the `PTL_LE_USE_ONCE` option is not inserted into the priority list. If a persistent list entry is posted to the priority list, it may cause multiple overflow events to be generated, one for every matching entry in the unexpected list. No permissions check is performed on a matching message in the unexpected list. No searching of the unexpected list is performed when a list entry is posted to the overflow list. When the list entry has been linked (inserted) into the specified list, a `PTL_EVENT_LINK` event is generated.

Discussion: Generally speaking, the user should attempt to insure that persistent list entries (or match list entries) are inserted before messages arrive that match them. Inserts of persistent entries could have unexpected performance and resource usage characteristics if a large unexpected list has accumulated, since a `PtlLEAppend()` that appends a persistent LE can cause multiple matches.

List Entry Type Constants (`ptl_list_t`)

<code>PTL_PRIORITY_LIST</code>	The priority list associated with a portal table entry
<code>PTL_OVERFLOW_LIST</code>	The overflow list associated with a portal table entry

Function Prototype for `PtlLEAppend`

```
int PtlLEAppend(ptl_handle_ni_t ni_handle,
               ptl_pt_index_t pt_index,
               const ptl_le_t *le,
               ptl_list_t ptl_list,
               void *user_ptr,
               ptl_handle_le_t *le_handle);
```

Arguments

<i>ni_handle</i>	input	The interface handle to use.
<i>pt_index</i>	input	The portal table index where the list entry should be appended.
<i>le</i>	input	Provides initial values for the user-visible parts of a list entry. Other than its use for initialization, there is no linkage between this structure and the list entry maintained by the API.
<i>ptl_list</i>	input	Determines whether the list entry is appended to the priority list or the overflow list.
<i>user_ptr</i>	input	A user-specified value that is associated with each command that can generate an event. The value does not need to be a pointer, but must fit in the space used by a pointer. This value (along with other values) is recorded in full events associated with operations on this list entry.
<i>le_handle</i>	output	On successful return, this location will hold the newly created list entry handle.

Return Codes

PTL_OK	Indicates success.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_NO_SPACE	Indicates that there is insufficient memory to allocate the match list entry.
PTL_LIST_TOO_LONG	Indicates that the resulting list is too long. The maximum length for a list is defined by the interface.

Discussion: *Tying commands to a user-defined value is useful at the target when the command needs to be associated with a data structure maintained by the process outside of the portals library. For example, an MPI implementation can set the *user_ptr* argument to the value of an MPI Request. This direct association allows for processing of list entries by the MPI implementation without a table look up or a search for the appropriate MPI Request.*

3.11.3. PtILEUnlink

The **PtILEUnlink()** function can be used to unlink a list entry from a list. If **PtILEUnlink()** returned **PTL_OK**, it is an error to use the list entry handle after the call to **PtILEUnlink()**. **PtILEUnlink()** will return **PTL_IN_USE** if the list entry is on the overflow list and has associated unexpected headers.

PtILEUnlink() is frequently used to implement the cancel of receive operations in higher level protocols. If the list entry handle passed to **PtILEUnlink()** has pending operations, e.g., an unfinished *put* operation or the list entry is in the overflow list and there are unexpected headers associated with the list entry, then **PtILEUnlink()** will return **PTL_IN_USE**, and the list entry will not be unlinked. An implementation must ensure that list entry handles remain valid for calls to **PtILEUnlink()** until the next call to **PtILEAppend()** after the last event associated with the list entry is delivered to an event queue or counting event. If the list entry has been unlinked before a call to **PtILEUnlink()** but before the next call to **PtILEAppend()**, **PtILEUnlink()** must return **PTL_IN_USE**.

**IMPLEMENTATION
NOTE 10:**

PtILEUnlink() and unlinked handles

PtILEUnlink() may be used to unlink list entries which are use-once. In this case, there is a race condition between a network operation causing a list entry to unlink and the list entry being explicitly unlinked. Requiring the handle to remain valid until the next call to **PtILEAppend()** allows higher level protocols to implement the serialization necessary to prevent such race conditions from impacting correctness. A Portals implementation does not need to limit the lifespan of handles to that specified. For example, a generation counter embedded in the handle may allow the handle to remain valid for the purposes of **PtILEUnlink()** for significantly longer than specified.

Function Prototype for PtILEUnlink

```
int PtILEUnlink(ptl_handle_le_t le_handle);
```

Arguments

le_handle **input** The list entry handle to be unlinked.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_IN_USE	Indicates that the list entry has pending operations and cannot be unlinked.

3.11.4. PtILESearch

The **PtILESearch()** function is used to search for a message in the unexpected list associated with a specific portal table entry specified by *pt_index* for the portal table for *ni_handle*. **PtILESearch()** uses the exact same search of the unexpected list as **PtILEAppend()**; however, the list entry specified in the **PtILESearch()** call is never linked into a priority list or an overflow list.

The **PtILESearch()** function can be called in two modes. If *ptl_search_op* is set to `PTL_SEARCH_ONLY`, the unexpected list is searched, and matching entries are left in the unexpected list. If *ptl_search_op* is set to `PTL_SEARCH_DELETE`, the unexpected list is searched, and matching items are deleted from the unexpected list. When used with `PTL_SEARCH_ONLY`, a `PTL_EVENT_SEARCH` event with *ni_fail_type* `PTL_NI_OK` is generated when a matching message is found in the unexpected list. When used with `PTL_SEARCH_DELETE`, the event that is generated corresponds to the type of operation that is found (e.g., `PTL_EVENT_PUT_OVERFLOW`, `PTL_EVENT_GET_OVERFLOW`, `PTL_EVENT_ATOMIC_OVERFLOW`, or `PTL_EVENT_FETCH_ATOMIC_OVERFLOW`).

Searches using LEs with the `PTL_LE_USE_ONCE` option set will cause at most one match, while searches without the `PTL_LE_USE_ONCE` option set (i.e., the LE is persistent) can cause multiple matches. When the `PTL_LE_USE_ONCE`

option is set, if a match is found, an event is generated as described above, and the search terminates; if no match is found, a `PTL_EVENT_SEARCH` event is generated with a failure indication of `PTL_NI_NO_MATCH`. When the `PTL_LE_USE_ONCE` option is not set, for each match found, an event is generated as previously described, and, at the end of the search, a `PTL_EVENT_SEARCH` event is generated with a failure indication of `PTL_NI_NO_MATCH` to indicate there are no remaining list entries matching the requested search in the unexpected list.

If the LE used for searching includes a counting event handle, and the `PTL_LE_USE_ONCE` option is set, the count of succeeding events (cf. Section 3.14.1) is incremented by one for the first match found on the unexpected list, and the search terminates. If no matches are found, the count of failing events is incremented by one. Consequently, when the `PTL_LE_USE_ONCE` option is set, incrementing the count of failing events indicates both that the search has completed and no events were found. If the LE used for searching includes a counting event handle, and the `PTL_LE_USE_ONCE` option is *not* set, the count of succeeding events is incremented by one for each match found in the unexpected list, and the count of failed events is incremented by one after all of the unexpected headers have been checked. Consequently, when the `PTL_LE_USE_ONCE` option is not set, incrementing the count of failing events indicates only that the search has completed; the count of succeeding events must be consulted to determine the number of matches found, if any.

No permissions check is performed during search; only matching criteria are used to determine if an event should be generated. Users should use the generated event data to perform any required permissions check.

Event generation for the search functions works just as it would for an append function. If a search is performed with full events disabled (either through option or through the absence of an event queue on the portal table entry), the search will succeed, but no full events will be generated. Status registers, however, are handled slightly differently for a search in that a `PtILESearch()` never causes a status register to be incremented.

Discussion: *Searches with persistent LEs could have unexpected performance and resource usage characteristics if a large unexpected list has accumulated, since a `PtILESearch()` that uses a persistent LE can cause multiple matches.*

List Entry Search Operation Constants (`ptl_search_op_t`)

<code>PTL_SEARCH_ONLY</code>	Use the LE/ME to search the unexpected list, without consuming an item in the list.
<code>PTL_SEARCH_DELETE</code>	Use the LE/ME to search the unexpected list and delete the item from the list.

Function Prototype for `PtILESearch`

```
int PtILESearch(ptl_handle_ni_t ni_handle,
               ptl_pt_index_t pt_index,
               const ptl_le_t *le,
               ptl_search_op_t ptl_search_op,
               void *user_ptr);
```

Arguments

<i>ni_handle</i>	input	The interface handle to use.
<i>pt_index</i>	input	The portal table index that should be searched.
<i>le</i>	input	Provides values for the user-visible parts of a list entry to use for searching.

<i>ptl_search_op</i>	input	Determines whether the function only searches the list or searches the list and deletes the matching entries from the list.
<i>user_ptr</i>	input	A user-specified value that is associated with each command that can generate an event. The value does not need to be a pointer, but must fit in the space used by a pointer. This value (along with other values) is recorded in full events associated with operations on this list entry.

Return Codes

PTL_OK	Indicates success.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.

3.12. Match List Entries and Matching Lists

Matching list entries add matching semantics to the basic list constructs. Each match list entry (ME) adds a set of match criteria to the basic memory region description in the list entry. The match criteria added can be used to reject incoming requests based on process identifier or the match bits provided in the request. The **PtIMEAppend()** function appends a single match list entry to the specified portal index and returns the match list entry handle. Matching list entries can be removed from a list using the **PtIMEUnlink()** function.

Like a list entry, a match list entry describes a memory region using a base address and length. A zero-length list entry may be created by setting *start* to NULL and *length* to 0. Zero-length buffers (NULL ME) are useful to record events. If truncation is not disabled, messages that are outside the bounds of the ME are truncated to zero bytes (e.g., zero-length buffers or an offset beyond the length of the ME). If the interface set the `PTL_TARGET_BIND_INACCESSIBLE` bit in the *features* field of the *actual* limits (See Section 3.6.1), then it is not a requirement for all of the memory described by the match list entry to be allocated or accessible within the application. For example, an application could create a match list entry that covers the entire virtual address range by setting *start* to NULL and *length* to `PTL_SIZE_MAX`, even though the entire region is not currently allocated (See Implementation Note 9). If an incoming operation (e.g., *put*) attempts to access an unallocated region of the ME, the implementation may either cause a segmentation fault of the application or may simply fail the operation. If a full event is delivered, it must set *ni_fail_type* to `PTL_NI_SEGV`. The target may, however, set the `PTL_ME_IS_ACCESSIBLE` option to indicate that the entire memory space described by the ME is accessible. If the match list entry sets the `PTL_IOVEC` option, the memory region(s) described by the **ptl_iovec_t** must all be accessible within the application. Note that a message with a length of zero bytes is not considered to have accessed memory. A message with an invalid address/offset pair that is of zero bytes will not cause a `PTL_NI_SEGV` *ni_fail_type* to be set in a full event or increment the number of failures in a **ptl_ct_event_t** structure.

Matching list entries can be appended to either the priority list or the overflow list associated with a portal table entry; however, when attached to an overflow list, additional semantics are implied that require the implementation to track messages that arrive in match list entries. Essentially, the memory region identified is provided to the implementation for use in managing unexpected messages; however, the application may use the match bits and other matching criteria to further constrain how these buffers are used. Buffers provided in the overflow list will post a full event (`PTL_EVENT_AUTO_UNLINK`) when the buffer space has been consumed, to notify the application that more buffer space may be needed. When the application is free to reuse the buffer (i.e. the implementation is done with it), another full event (`PTL_EVENT_AUTO_FREE`) will be posted. The `PTL_EVENT_AUTO_FREE` full event will be posted after *all* other events associated with the buffer have been posted to the event queue.

Match list entries provide semantics to allow the target of an operation to determine data placement, through the use of locally managed offsets and the ability to automatically unlink when a minimum free space value is reached. These semantics are useful for the efficient implementation of unexpected or active message semantics. The locally managed offset is changed when data is delivered into a given match list entry and, likewise, the minimum free test is applied when data is delivered into a match list entry. Therefore, unless the `PTL_ME_LOCAL_INC_UH_RLENGTH` option is set, when a locally managed, persistent match list entry is appended to the priority list and matches headers in the unexpected headers list, the match list's local offset will not move based on the unexpected headers.

In the case where a match list operation fails (e.g., a hardware error), there is the possibility that the local offset is incremented without data being delivered. If an error occurred and the locally managed offset was incremented, the min free test must still be performed.

Incoming match bits are compared to the match bits stored in the match list entry using the ignore bits as a mask. An optimized version of this is shown in the following code fragment:

```
((incoming_bits ^ match_bits) & ~ignore_bits) == 0
```

Discussion: *It is the responsibility of the application to ensure that the implementation has sufficient buffer space to manage unexpected messages. Failure to do will cause messages to be dropped. The `PTL_EVENT_ACK` at the initiator will indicate the failure as described in Section 3.13.3. Note that overflow events can readily exhaust the event queue. Proper use of the API will generally require the application to post at least two (and typically several) buffers so that the application has time to notice the `PTL_EVENT_AUTO_UNLINK` and replace the buffer.*

It is the responsibility of the implementation to determine when a buffer unlinked from an overflow list can be reused. It must note that it is no longer holding state associated with the buffer and deliver a `PTL_EVENT_AUTO_FREE` full event after all other events associated with that buffer have been delivered.

Match list entries may only be appended to a matching network interface. The interpretation of the `match_id` field in a match list entry is determined by whether the network interface is physically or logically addressed.

3.12.1. The Match List Entry Type

The `ptl_me_t` type defines the visible parts of a match list entry. Values of this type are used to initialize and update the match list entries.

```
typedef struct {
    void *start;
    ptl_size_t length;
    ptl_handle_ct_t ct_handle;
    ptl_uid_t uid;
    unsigned int options;
    ptl_process_t match_id;
    ptl_match_bits_t match_bits;
    ptl_match_bits_t ignore_bits;
    ptl_size_t min_free;
} ptl_me_t;
```

Members

<i>start, length</i>	Specify the memory region associated with the match list entry. The <i>start</i> member specifies the starting address for the memory region and the <i>length</i> member specifies the length of the region. There are no restrictions on buffer alignment, the starting address or the length of the region; although messages that are not natively aligned (e.g. to a four byte or eight byte boundary) may be slower (i.e., lower bandwidth and/or longer latency) on some implementations.
<i>ct_handle</i>	A handle for counting events associated with the memory region. If this argument is <code>PTL_CT_NONE</code> , operations performed on this match list entry are not counted.
<i>min_free</i>	When the unused portion of a match list entry (<code>length - local offset</code>) falls below this value, the match list entry automatically unlinks. A <i>min_free</i> value of 0 disables the <i>min_free</i> capability (the free space cannot fall below 0). This value is only used if <code>PTL_ME_MANAGE_LOCAL</code> is set.
<i>uid</i>	Specifies the usage ID that may access this match list entry. The usage ID may be set to a wildcard (<code>PTL_UID_ANY</code>). If the access control check fails, then the message is dropped without modifying Portals state. This is treated as a permissions failure and the status register indexed by <code>PTL_SR_PERMISSION_VIOLATIONS</code> is incremented. This failure is also indicated to the initiator. If a full event is delivered to the initiator, the <i>ni_fail_type</i> in the <code>PTL_EVENT_ACK</code> full event must be set to <code>PTL_NI_PERM_VIOLATION</code> .
<i>options</i>	Specifies the behavior of the match list entry. The following options can be selected: enable <i>put</i> operations (yes or no), enable <i>get</i> operations (yes or no), offset management (local or remote), message truncation (yes or no), acknowledgment (yes or no), use scatter/gather vectors and control event delivery. Values for this argument can be constructed using a bitwise OR of the following values:
<code>PTL_ME_OP_PUT</code>	Specifies that the match list entry will respond to <i>put</i> operations. By default, match list entries reject <i>put</i> operations. If a <i>put</i> operation targets a list entry where <code>PTL_ME_OP_PUT</code> is not set, it is treated as an operations failure and <code>PTL_SR_OPERATION_VIOLATIONS</code> is incremented. If a full event is delivered to the initiator, the <i>ni_fail_type</i> in the <code>PTL_EVENT_ACK</code> event must be set to <code>PTL_NI_OP_VIOLATION</code> .
<code>PTL_ME_OP_GET</code>	Specifies that the match list entry will respond to <i>get</i> operations. By default, match list entries reject <i>get</i> operations. If a <i>get</i> operation targets a list entry where <code>PTL_ME_OP_GET</code> is not set, it is treated as an operations failure and <code>PTL_SR_OPERATION_VIOLATIONS</code> is incremented. If a full event is delivered to the initiator, the <i>ni_fail_type</i> in the <code>PTL_EVENT_ACK</code> event must be set to <code>PTL_NI_OP_VIOLATION</code> .

Note: It is not considered an error to have a match list entry that responds to both *put* and *get* operations. In fact, a match list entry must be configured to respond to both *put* and *get* operations to properly handle a `PtlFetchAtomic()` or `PtlSwap()` operation.

PTL_ME_MANAGE_LOCAL	<p>Specifies that the offset used in accessing the memory region is managed locally. By default, the offset is in the incoming message. When the offset is maintained locally, the offset is incremented by the length of the request so that the next operation (<i>put</i> and/or <i>get</i>) will access the next part of the memory region.</p> <p>Note that only one offset variable exists per match list entry. If both <i>put</i> and <i>get</i> operations are performed on a match list entry, the value of that single variable is updated each time.</p>
PTL_ME_LOCAL_INC_UH_RLENGTH	<p>Specifies that the local offset should be incremented by the requested length of every matching header found in the unexpected headers list. When this option is set and a locally managed, persistent match list entry is appended to the priority list, the local offset is incremented by the requested length of every matching header found in the unexpected headers list as long as the <i>min_free</i> condition is respected. Note that the <i>min_free</i> semantics is checked for each matching header found in the unexpected headers list. If, after moving the local offset, the <i>min_free</i> condition is no longer valid, the match list entry is not appended to the priority list and a PTL_EVENT_AUTO_UNLINK event is generated. This option is only relevant when the PTL_ME_MANAGE_LOCAL option is set. When the PTL_ME_LOCAL_INC_UH_RLENGTH option is set, the PTL_ME_MAY_ALIGN option is ignored.</p>
PTL_ME_NO_TRUNCATE	<p>Specifies that the length provided in the incoming request cannot be reduced to match the memory available in the region. This will cause the matching to fail for a match list entry and continue with the next entry. (The memory available in a memory region is determined by subtracting the offset from the length of the memory region.) Note that zero length messages are not checked for truncation, meaning that they will never fail to match due to an offset that does not pass the truncation check. By default, if the length in the incoming operation is greater than the amount of memory available, the operation is truncated.</p>
PTL_ME_USE_ONCE	<p>Specifies that the match list entry will only be used once and then automatically unlinked by the implementation. If this option is not set, the match list entry persists until it is explicitly unlinked or another unlink condition is triggered.</p>
PTL_ME_MAY_ALIGN	<p>Indicate that messages deposited into this match list entry may be aligned by the implementation to a performance optimizing boundary. Essentially, this is a performance hint to the implementation to indicate that the application does not care about the specific placement of the data. This option is only relevant when the PTL_ME_MANAGE_LOCAL option is set. The PTL_ME_MAY_ALIGN option is ignored if the PTL_ME_LOCAL_INC_UH_RLENGTH option is also set.</p>
PTL_ME_UNEXPECTED_HDR_DISABLE	<p>Specifies that the header for a message delivered to this match list entry should not be added to the unexpected list. This option only has meaning if the match list entry is inserted into the overflow list. By creating a match list entry which truncates messages to zero bytes, disables comm events, and sets this option, a user may create a match list entry which consumes no target side resources. A list entry with this flag set does not generate PTL_EVENT_AUTO_FREE events.</p>

PTL_IOVEC	Specifies that the <code>start</code> argument is a pointer to an array of type <code>ptl_iovec_t</code> (Section 3.10.2) and the <code>length</code> argument is the length of the array. This allows for a scatter/gather capability for match list entries. A scatter/gather match list entry behaves exactly as a match list entry that describes a single virtually contiguous region of memory. All other semantics are identical.
PTL_ME_IS_ACCESSIBLE	Indicate that this match list entry only contains memory addresses that are accessible by the application.
PTL_ME_EVENT_LINK_DISABLE	Specifies that this match list entry should not generate a <code>PTL_EVENT_LINK</code> full event indicating the list entry successfully linked.
PTL_ME_EVENT_COMM_DISABLE	Specifies that this match list entry should not generate full events that indicate a communication operation. This includes <code>PTL_EVENT_GET</code> , <code>PTL_EVENT_PUT</code> , <code>PTL_EVENT_ATOMIC</code> , <code>PTL_EVENT_FETCH_ATOMIC</code> , and <code>PTL_EVENT_SEARCH</code> .
PTL_ME_EVENT_FLOWCTRL_DISABLE	Specifies that this match list entry should not generate a <code>PTL_EVENT_PT_DISABLED</code> full event that indicate a flow control failure.
PTL_ME_EVENT_SUCCESS_DISABLE	Specifies that this match list entry should not generate full events if the <i>ni_fail_type</i> would be <code>PTL_OK</code> . This flag does not affect counting events. Disabling full events for successful operations is useful in scenarios when a counting event is sufficient for completion, but more information is needed for error recovery.
PTL_ME_EVENT_OVER_DISABLE	Specifies that this match list entry should not generate overflow list full events. This includes <code>PTL_EVENT_PUT_OVERFLOW</code> , <code>PTL_EVENT_GET_OVERFLOW</code> , <code>PTL_EVENT_ATOMIC_OVERFLOW</code> , and <code>PTL_EVENT_FETCH_ATOMIC_OVERFLOW</code> .
PTL_ME_EVENT_UNLINK_DISABLE	Specifies that this match list entry should not generate auto-unlink (<code>PTL_EVENT_AUTO_UNLINK</code>) or free (<code>PTL_EVENT_AUTO_FREE</code>) full events.
PTL_ME_EVENT_CT_COMM	Enable the counting of communication events (<code>PTL_EVENT_PUT</code> , <code>PTL_EVENT_GET</code> , <code>PTL_EVENT_ATOMIC</code> , <code>PTL_EVENT_FETCH_ATOMIC</code> and <code>PTL_EVENT_SEARCH</code>).
PTL_ME_EVENT_CT_OVERFLOW	Enable the counting of overflow events (<code>PTL_EVENT_PUT_OVERFLOW</code> , <code>PTL_EVENT_GET_OVERFLOW</code> , <code>PTL_EVENT_ATOMIC_OVERFLOW</code> , <code>PTL_EVENT_FETCH_ATOMIC_OVERFLOW</code>).
PTL_ME_EVENT_CT_BYTES	By default, counting events count events. When set, this option causes bytes to be counted instead for success events. Byte counts must be incremented exactly once per operation. The increment is by the number of bytes counted (<i>mlength</i>). Failure events always increment the count by one.
<i>match_id</i>	Specifies the match criteria for the process identifier of the requester. The constants <code>PTL_PID_ANY</code> and <code>PTL_NID_ANY</code> can be used to wildcard either of the physical identifiers in the <code>ptl_process_t</code> structure, or <code>PTL_RANK_ANY</code> can be used to wildcard the rank for logical addressing.
<i>match_bits, ignore_bits</i>	Specify the match criteria to apply to the match bits in the incoming request. The <i>ignore_bits</i> are used to mask out insignificant bits in the incoming match bits. The resulting bits are then compared to the match list entry's match bits to determine if the incoming request meets the match criteria.

Discussion: *The default behavior from Portals 3.3 (no truncation and locally managed offsets) has been changed to match the default semantics of the list entry, which does not provide matching.*

When the `PTL_ME_USE_ONCE` option is set, an event associated with a target side operation (e.g. a `PTL_EVENT_PUT` event) also implies that the associated match list entry has unlinked; hence, it is safe on these match list entries to set the `PTL_ME_EVENT_UNLINK_DISABLE` option.

`PTL_ME_EVENT_FLOWCTRL_DISABLE` only disables flow control events which are the direct result of an incoming message matching the current match list entry. This includes a message matching the match list entry but the associated event queue is full or a message matching a match list entry in the overflow list but the unexpected headers list is full. If flow control is enabled on the portal table entry and a message does not match in either the priority or overflow lists, a `PTL_EVENT_PT_DISABLED` event is always generated.

Although the `MD`, `ME`, and `LE` can all map inaccessible memory, only the `ME` and `LE` have an option to allow the user to indicate to the implementation that the entire region is accessible. This is because the typical usage model for the `MD` is expected to bind inaccessible memory, while a very common usage model for both the `ME` and `LE` is expected to only use accessible memory.

3.12.2. PtlMEAppend

The `PtlMEAppend()` function creates a single match list entry. If `PTL_PRIORITY_LIST` or `PTL_OVERFLOW_LIST` is specified by `ptl_list`, this entry is appended to the end of the appropriate list specified by `ptl_list` associated with the portal table entry specified by `pt_index` for the portal table for `ni_handle`.

When a match list entry is posted to the priority list, the unexpected list is searched to see if a matching message has been delivered in the overflow list prior to the posting of the match list entry. If so, an appropriate overflow event is generated, the matching header is removed from the unexpected list, and a match list entry with the `PTL_ME_USE_ONCE` option is not inserted into the priority list. If a persistent match list entry is posted to the priority list, it may cause multiple overflow events to be generated, one for every matching entry in the unexpected list. If a persistent, locally managed match list entry with option `PTL_ME_LOCAL_INC_UH_RLENGTH` is posted to the priority list, prior to the posting of the match list entry, the local offset will be incremented by the requested length of each matched header found in the unexpected headers list as long as the `min_free` condition is respected. In this case, if the `PTL_ME_MAY_ALIGN` option is also set, the `PTL_ME_MAY_ALIGN` option is ignored. In all cases, no permissions checking is performed on a matching message in the unexpected list. No searching of the unexpected list is performed when a match list entry is posted to the overflow list. When the list entry has been linked (inserted) into the specified list, a `PTL_EVENT_LINK` event is generated.

Discussion: Generally speaking, the user should attempt to insure that persistent match list entries (or simple list entries) are inserted before messages arrive that match them. Appending of persistent entries could have unexpected performance and resource usage characteristics if a large unexpected list has accumulated, since a `PtlMEAppend()` that appends a persistent `ME` can cause multiple matches.

See the `PtlLEAppend()` definition in Section 3.11.2 for the definition of `ptl_list_t`.

Function Prototype for PtlMEAppend

```
int PtlMEAppend(ptl_handle_ni_t ni_handle,
                ptl_pt_index_t pt_index,
                const ptl_me_t *me,
                ptl_list_t ptl_list,
                void *user_ptr,
                ptl_handle_me_t *me_handle);
```

Arguments

<i>ni_handle</i>	input	The interface handle to use.
<i>pt_index</i>	input	The portal table index where the match list entry should be appended.
<i>me</i>	input	Provides initial values for the user-visible parts of a match list entry. Other than its use for initialization, there is no linkage between this structure and the match list entry maintained by the API.
<i>ptl_list</i>	input	Determines whether the match list entry is appended to the priority list or the overflow list.
<i>user_ptr</i>	input	A user-specified value that is associated with each command that can generate an event. The value does not need to be a pointer, but must fit in the space used by a pointer. This value (along with other values) is recorded in full events associated with operations on this match list entry.
<i>me_handle</i>	output	On successful return, this location will hold the newly created match list entry handle.

Return Codes

PTL_OK	Indicates success.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_NO_SPACE	Indicates that there is insufficient memory to allocate the match list entry.
PTL_LIST_TOO_LONG	Indicates that the resulting list is too long. The maximum length for a list is defined by the interface.

Discussion: *Tying commands to a user-defined value is useful at the target when the command needs to be associated with a data structure maintained by the process outside of the portals library. For example, an MPI implementation can set the *user_ptr* argument to the value of an MPI Request. This direct association allows for processing of match list entries by the MPI implementation without a table look up or a search for the appropriate MPI Request.*

IMPLEMENTATION NOTE 11:

Checking *match_id* Argument

Checking whether a *match_id* is a valid process identifier may require global knowledge. However, **PtIMEAppend()** is not meant to cause any communication with other nodes in the system. Therefore, **PTL_ARG_INVALID** may not be returned in some cases where it would seem appropriate.

3.12.3. PtIMEUnlink

The **PtIMEUnlink()** function can be used to unlink a match list entry from a list. If **PtIMEUnlink()** returned **PTL_OK**, it is an error to use the match list entry handle after the call to **PtIMEUnlink()**. **PtIMEUnlink()** should return **PTL_IN_USE** if the match list entry is on the overflow list and has associated unexpected headers.

If the ME used for searching includes a counting event handle, and the `PTL_ME_USE_ONCE` option is set, the count of succeeding events (cf. Section 3.14.1) is incremented by one for the first match found on the unexpected list, and the search terminates. If no matches are found, the count of failing events is incremented by one. Consequently, when the `PTL_ME_USE_ONCE` option is set, incrementing the count of failing events indicates both that the search has completed and no events were found. If the ME used for searching includes a counting event handle, and the `PTL_ME_USE_ONCE` option is *not* set, the count of succeeding events is incremented by one for each match found in the unexpected list, and the count of failed events is incremented by one after all of the unexpected headers have been checked. Consequently, when the `PTL_ME_USE_ONCE` option is not set, incrementing the count of failing events indicates only that the search has completed; the count of succeeding events must be consulted to determine the number of matches found, if any.

No permissions checking is performed during search; only matching criteria are used to determine if an event should be generated. Users should use the generated event data to perform any required permissions check.

Event generation for the search functions works just as it would for an append function. If a search is performed with full events disabled (either through option or through the absence of an event queue on the portal table entry), the search will succeed, but no events will be generated. Status registers, however, are handled slightly differently for a search in that a `PtIMESearch()` never causes a status register to be incremented.

See the `PtILESearch()` definition in Section 3.11.4 for the definition of `ptl_search_op` and important notes associated with implementing and using `PtIMESearch()`.

Function Prototype for PtIMESearch

```
int PtIMESearch(ptl_handle_t ni_handle,
               ptl_pt_index_t pt_index,
               const ptl_me_t *me,
               ptl_search_op_t ptl_search_op,
               void *user_ptr);
```

Arguments

<i>ni_handle</i>	input	The interface handle to use.
<i>pt_index</i>	input	The portal table index that should be searched.
<i>me</i>	input	Provides values for the user-visible parts of a match list entry to use for searching.
<i>ptl_search_op</i>	input	Determines whether the function only searches the list or searches the list and deletes the matching entries from the list.
<i>user_ptr</i>	input	A user-specified value that is associated with each command that can generate an event. The value does not need to be a pointer, but must fit in the space used by a pointer. This value (along with other values) is recorded in full events associated with operations on this match list entry.

Return Codes

PTL_OK	Indicates success.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.

3.13. Events and Event Queues

Event queues are used to log operations performed on memory descriptors, list entries, match list entries, or portal table entries. In particular, they signal the end of a data transmission into or out of a memory region. They can also be used to hold acknowledgments for completed *put* operations and indicate when a list entry has been unlinked. Multiple memory descriptors or portal table entries can share a single event queue.

In addition to the `ptl_handle_eq_t` type, the Portals API defines two types associated with full events: The `ptl_event_kind_t` type is an integral type which defines the kinds of events that can be stored in an event queue. The `ptl_event_t` type defines the structure that is placed into event queues.

The Portals API provides five functions for dealing with event queues: The `PtlEQAlloc()` function is used to allocate the API resources needed for an event queue, the `PtlEQFree()` function is used to release these resources, the `PtlEQGet()` function can be used to get the next full event from an event queue, the `PtlEQWait()` function can be used to block a process (or thread) until an event queue has at least one full event, and the `PtlEQPoll()` function can be used to test or wait on multiple event queues, which may be associated with different logical network interfaces if they all belong to a single physical network interface. Thread safety requires that an event can only be returned to one `PtlEQGet()`, `PtlEQPoll()`, or `PtlEQWait()` call.

3.13.1. Kinds of Events

The Portals API defines sixteen types of events that can be logged:

Event Type Constants (`ptl_event_kind_t`)

<code>PTL_EVENT_GET</code>	A <i>get</i> operation completed at the <i>target</i> . Portals will not read from memory on behalf of this operation once this event has been logged.
<code>PTL_EVENT_GET_OVERFLOW</code>	A list entry posted by <code>PtlEAppend()</code> or <code>PtlMEAppend()</code> matched a <i>get</i> header in the unexpected list.
<code>PTL_EVENT_PUT</code>	A <i>put</i> operation completed at the <i>target</i> . Portals will not alter memory on behalf of this operation once this event has been logged.
<code>PTL_EVENT_PUT_OVERFLOW</code>	A list entry posted by <code>PtlEAppend()</code> or <code>PtlMEAppend()</code> matched a <i>put</i> header in the unexpected list.
<code>PTL_EVENT_ATOMIC</code>	An <i>atomic</i> operation that does not return data to the <i>initiator</i> completed at the <i>target</i> . Portals will not read from or alter memory on behalf of this operation once this event has been logged.
<code>PTL_EVENT_ATOMIC_OVERFLOW</code>	A list entry posted by <code>PtlEAppend()</code> or <code>PtlMEAppend()</code> matched an <i>atomic</i> header in the unexpected list for an operation which does not return data to the <i>initiator</i> .
<code>PTL_EVENT_FETCH_ATOMIC</code>	An <i>atomic</i> operation that returns data to the initiator completed at the <i>target</i> . These include <code>PtlFetchAtomic()</code> and <code>PtlSwap()</code> . Portals will not read from or alter memory on behalf of this operation once this event has been logged.
<code>PTL_EVENT_FETCH_ATOMIC_OVERFLOW</code>	A list entry posted by <code>PtlEAppend()</code> or <code>PtlMEAppend()</code> matched an <i>atomic</i> header in the unexpected list for an operation which returns data to the <i>initiator</i> .

PTL_EVENT_REPLY	A <i>reply</i> operation has completed at the <i>initiator</i> , either due to a <i>get</i> operation or an <i>atomic</i> which returned data to the initiator. This event is logged after the data (if any) from the reply has been written into the memory descriptor. Receipt of a PTL_EVENT_REPLY indicates remote completion of the operation.
PTL_EVENT_SEND	A <i>put</i> or <i>atomic</i> has completed at the <i>initiator</i> . This event is logged after it is safe to reuse the buffer, but does not mean the message has been processed by the <i>target</i> .
PTL_EVENT_ACK	An <i>acknowledgment</i> was received. This event is logged when the acknowledgment is received. Receipt of a PTL_EVENT_ACK indicates remote completion of the operation. Remote completion indicates that local completion has also occurred.
PTL_EVENT_PT_DISABLED	Resources exhaustion has occurred on this portal table entry, which has entered a flow control situation. See Section 2.8.
PTL_EVENT_LINK	A list entry posted by PtILEAppend() or PtIMEAppend() has successfully linked into the specified list.
PTL_EVENT_AUTO_UNLINK	A list entry/match list entry was automatically unlinked (Sections 3.12.2 and 3.11.2). A PTL_EVENT_AUTO_UNLINK event is generated even if the list entry/match list entry passed into the PtILEAppend()/PtIMEAppend() operation was marked with the PTL_LE_USE_ONCE/PTL_ME_USE_ONCE option and found a corresponding unexpected message before being “linked” into the priority list. A PTL_EVENT_AUTO_UNLINK must be delivered after all PTL_EVENT_GET, PTL_EVENT_PUT, PTL_EVENT_ATOMIC, and PTL_EVENT_FETCH_ATOMIC events associated with the list entry/match list entry have been delivered.
PTL_EVENT_AUTO_FREE	A list entry/match list entry previously automatically unlinked from the overflow list is now free to be reused by the application. A PTL_EVENT_AUTO_FREE event is generated when Portals will not generate any further events which resulted from messages delivered into the specified overflow list entry. This also indicates that the unexpected list contains no more items associated with this entry. A list entry/match list entry which disabled unexpected headers will not generate this event, even if placed in the overflow list.
PTL_EVENT_SEARCH	A PtILESearch() or PtIMESearch() call completed. If a matching message was found in the overflow list, PTL_NI_OK is returned in the <i>ni_fail_type</i> field of the event and the event queue entries are filled in as if it were an overflow event. If no matching message is found, or if matching messages were found but no more list entries would match the requested search (i.e., the PtIMESearch() call is persistent and generated an event for all matches on the list), a failure is recorded in the <i>ni_fail_type</i> field using PTL_NI_NO_MATCH, the <i>user_ptr</i> is filled in correctly, and the other fields are undefined. In cases where matches were found and events generated, this indicates to the user that all matching list entries have been reported.
PTL_EVENT_ERROR	An error occurred that is not specified or cannot return all of the required fields in a valid error type.

Discussion: *PTL_EVENT_ERROR* is intended to be used in cases where unspecified errors may be detectable and recoverable by the application. For example, file systems may be able to recover from errors that cannot be fully described by a Portals implementation.

Overflow events are used to indicate that a message matching the list entry or match list entry posted by **PtILEAppend()** or **PtIMEAppend()** was previously delivered into the overflow list and its header was found in the unexpected list (See Section 2.5). The operation was processed as specified by the list entry in the overflow list to which it matched, meaning that all, some, or none of the message may have been written to or read from the matching list entry in the overflow list. The full event's *start* will point to the start of the message (or where the message was read, in the case of a *get* operation). The *rlength* and *mlength* of the full event may be used to determine whether the message was fully delivered or truncated.

Discussion: When an application wishes to record unexpected messages, it may place an entry on the overflow list which has no memory associated with it and truncates all messages to zero bytes. The *hdr_data* field, along with a higher-level protocol, may be used to complete the transaction at a later time. In the case of MPI, a number of match list entries on the overflow list with locally managed offsets may additionally be used to optimize unexpected short messages.

3.13.2. Event Occurrence

The diagrams in Figure 3-1 show when events occur in relation to portals operations and whether they are recorded on the *initiator* or the *target* side. Note that local and remote events are not synchronized or ordered with respect to each other.

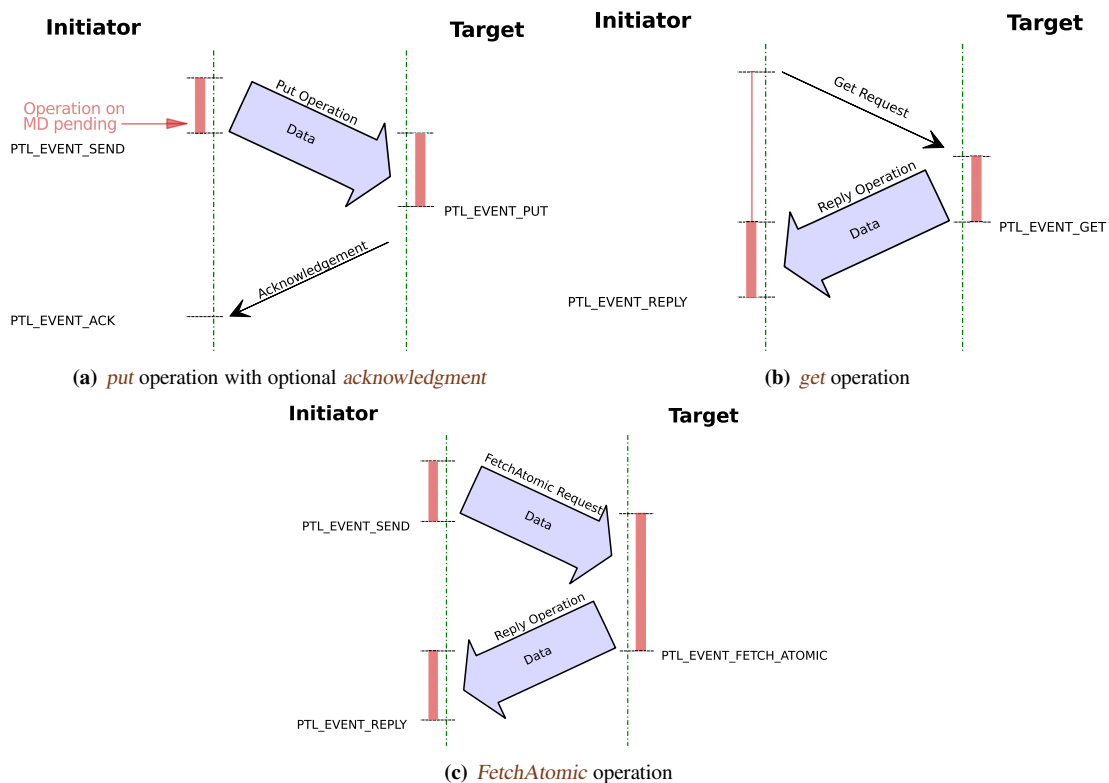


Figure 3-1. Portals Operations and Event Types: The red bars indicate the times a local memory descriptor is considered to be in use by the system; i.e., it has operations pending. Users should not modify memory descriptors or match list entries during those periods.

Figure 3-1(a) shows the events that are generated for a *put* operation including the optional *acknowledgment*. The diagram shows which events are generated at the *initiator* and the *target* side of the *put* operation. Figure 3-1(b) shows the corresponding events for a *get* operation, and Figure 3-1(c) shows the events generated for an *atomic* operation.

When the initiator of an operation receives a remote completion event (e.g. `PTL_EVENT_ACK`), local completion is also implied. While no ordering is required between local and remote completion events at the initiator (i.e. there is no guaranteed ordering between `PTL_EVENT_SEND` and `PTL_EVENT_ACK` for the same operation), a user may reuse a buffer after either the local or remote completion event is received.

If, as a result of any of the operations shown in the diagrams of Figure 3-1, a match list entry is unlinked, then a `PTL_EVENT_AUTO_UNLINK` event is generated on the *target*. This is not shown in the diagrams. No *initiator* events are generated if the memory descriptor does not have an attached event queue. Similarly, no *target* events are generated if the portal table entry associated with the matched list entry does not have an attached event queue. See the description of `PTL_EQ_NONE` on page 48 of Section 3.10.1) for more information. The various types of events can also be disabled by type (e.g. see the description of `PTL_ME_EVENT_COMM_DISABLE` and `PTL_ME_EVENT_UNLINK_DISABLE` on page 69, also in Section 3.12.1.).

Table 3-2 summarizes the portals event types and where each event type may be generated.

Table 3-2. Event Type Summary: A list of event types and where (*initiator* or *target*) they can occur.

Event Type	<i>initiator</i>	<i>target</i>
<code>PTL_EVENT_GET</code>		•
<code>PTL_EVENT_GET_OVERFLOW</code>		•
<code>PTL_EVENT_PUT</code>		•
<code>PTL_EVENT_PUT_OVERFLOW</code>		•
<code>PTL_EVENT_ATOMIC</code>		•
<code>PTL_EVENT_ATOMIC_OVERFLOW</code>		•
<code>PTL_EVENT_FETCH_ATOMIC</code>		•
<code>PTL_EVENT_FETCH_ATOMIC_OVERFLOW</code>		•
<code>PTL_EVENT_REPLY</code>	•	
<code>PTL_EVENT_SEND</code>	•	
<code>PTL_EVENT_ACK</code>	•	
<code>PTL_EVENT_PT_DISABLED</code>		•
<code>PTL_EVENT_LINK</code>		•
<code>PTL_EVENT_AUTO_UNLINK</code>		•
<code>PTL_EVENT_AUTO_FREE</code>		•
<code>PTL_EVENT_SEARCH</code>		•
<code>PTL_EVENT_ERROR</code>		•

3.13.3. Failure Notification

There are three ways in which operations may fail to complete successfully: the system (hardware or software) can fail in a way that makes the message undeliverable, a permissions violation can occur at the target, or resources can be exhausted at a target that has enabled flow-control. In any other scenario, every operation that is started will eventually complete. While an operation is in progress, the memory on the *target* associated with the operation should not be viewed (in the case of a *put* or a *reply*) or altered on the *initiator* side (in the case of a *put* or *get*). Operation completion, whether successful or unsuccessful, is final. That is, when an operation completes, the memory associated with the operation will no longer be read or altered by the operation. A network interface can use the integral type `ptl_ni_fail_t` to define specific information regarding the failure of the operation and record this information in the `ni_fail_type` field of a full event. Portals defines a number of event failure constants:

Event Failure Type Constants (`ptl_ni_fail_t`)

`PTL_NI_OK`

The operation causing the event was successful.

PTL_NI_UNDELIVERABLE	Indicates a system failure that prevents message delivery.
PTL_NI_PT_DISABLED	Indicates that the portal table entry at the <i>target</i> was disabled and did not process the operation, either because the entry was disabled with PtIPTDisable() or because the entry provides flow control and a resource has been exhausted. This failure type should only be returned on <i>initiator</i> events.
PTL_NI_DROPPED	Indicates that the message associated with this full event was dropped at the <i>target</i> for reasons other than a disabled portal table entry. This failure type should only be returned on <i>initiator</i> events.
PTL_NI_PERM_VIOLATION	Indicates that the remote Portals addressing has indicated a permissions violation for the operation that caused this event. This failure type should only be returned on <i>initiator</i> events.
PTL_NI_OP_VIOLATION	Indicates that the remote Portals addressing has indicated an operation violation for the operation that caused this event. This failure type should only be returned on <i>initiator</i> events.
PTL_NI_SEGV	A message attempted to access inaccessible memory.
PTL_NI_NO_MATCH	A search did not find an entry or all entries were found in the unexpected list and all corresponding events were generated.

To allow PTL_EVENT_SEND events to be local operations, all errors requiring remote information are delivered in PTL_EVENT_ACK or PTL_EVENT_REPLY events. This means that a PTL_EVENT_ACK will be delivered if it is requested, except when: 1) flow control is not enabled on the target portal table entry and the message does not match in either the priority list or overflow list or the message matches in the overflow list and the unexpected headers list is full, or 2) a locally generated failure is delivered in the PTL_EVENT_SEND. Certain classes of failures (e.g. a PTL_NI_UNDELIVERABLE that results from the network bifurcating) may require a local timeout to guarantee that the PTL_EVENT_ACK or PTL_EVENT_REPLY event is delivered.

Discussion: *Because remote errors are indicated in the PTL_EVENT_ACK or PTL_EVENT_REPLY events, the PTL_EVENT_SEND event only guarantees that the Portals implementation will not touch the buffer again. If the user intends to recover from a remote error, then the user cannot determine that an operation is done until the PTL_EVENT_ACK or PTL_EVENT_REPLY event is received.*

<p>IMPLEMENTATION NOTE 12:</p>	<p><u>Completion of portals operations</u></p> <p>Portals guarantees that every operation started will finish with an event if events are not disabled. While this document cannot enforce or recommend a suitable time, a quality implementation will keep the amount of time between an operation initiation and a corresponding event as short as possible. That includes operations that do not complete successfully. Timeouts of underlying protocols should be chosen accordingly.</p>
--	---

3.13.4. The Event Structure

An event queue contains **ptl_event_t** structures. An operation on the *target* needs information about the local match list entry modified, the initiator of the operation and the operation itself. The *initiator*, in contrast, can track all information about the attempted operation; however, it does need the result of the operation and a pointer to resolve back to the local structure tracking the information about the operation.

Many fields in the `ptl_event_t` structure only have meaning for a subset of the event types. Further, an implementation is not required to provide all fields in the `ptl_event_t` structure when the event is reporting an error. Table 3-3 defines which fields are defined in both success and error conditions.

```
typedef struct {
    void *start;
    void *user_ptr;
    ptl_hdr_data_t hdr_data;
    ptl_match_bits_t match_bits;
    ptl_size_t rlength;
    ptl_size_t mlength;
    ptl_size_t remote_offset;
    ptl_uid_t uid;
    ptl_process_t initiator; /* nid, pid or rank */
    ptl_event_kind_t type;
    ptl_list_t ptl_list;
    ptl_pt_index_t pt_index;
    ptl_ni_fail_t ni_fail_type;
    ptl_op_t atomic_operation;
    ptl_datatype_t atomic_type;
} ptl_event_t;
```

Members

<i>start</i>	<p>The starting location (virtual, byte address) where the message has been placed. The <i>start</i> variable is the sum of the <i>start</i> variable in the list entry and the offset used for the operation. The offset can be determined by the operation (Section 3.15) for a remote managed match list entry or by the local memory descriptor (Section 3.12). In the case of iovecs, the <i>start</i> is still the first address where the message was placed or read from, even if multiple iovec entries were used.</p> <p>When an append call matches a message that has arrived in the overflow list, the start address points to the address in the overflow list where the matching message resides. This may require the application to copy the message to the desired buffer.</p>
<i>user_ptr</i>	<p>The user-specified value associated with the local command that generated the full event. Note that, unlike <i>hdr_data</i>, the <i>user_ptr</i> is a locally-generated value. For example, the <i>user_ptr</i> for a full event of type PTL_EVENT_PUT is the <i>user_ptr</i> specified to the associated call to PtILEAppend() or PtIMEAppend(). For further discussion of <i>user_ptr</i>, see Section 3.12.2.</p>
<i>hdr_data</i>	64 bits of out-of-band user data (Section 3.15.2).
<i>match_bits</i>	The match bits specified by the <i>initiator</i> . This field should be set to 0 if the event is associated with a non-matching list entry.
<i>rlength</i>	The length (in bytes) specified in the request.
<i>mlength</i>	<p>The length (in bytes) of the data that was manipulated by the operation. For PTL_EVENT_SEND events, the manipulated length is the number of bytes sent, which may be larger than the number of bytes delivered (which can be determined by examining the <i>mlength</i> of the associated PTL_EVENT_ACK event). For PTL_EVENT_PUT, PTL_EVENT_GET, PTL_EVENT_ATOMIC, or PTL_EVENT_FETCH_ATOMIC events, the manipulated length is the number of bytes manipulated (delivered into or read from memory) at the target, which may be less than the <i>rlength</i> in the case of truncated operations. For PTL_EVENT_SEARCH and the overflow events, the manipulated length is the same value as the <i>mlength</i> returned in the corresponding PTL_EVENT_PUT, PTL_EVENT_GET, PTL_EVENT_ATOMIC, or PTL_EVENT_FETCH_ATOMIC event generated when the operation completed in the list entry on the overflow list.</p>
<i>remote_offset</i>	<p>The offset requested/used by the other end of the communication. At the initiator, this is the displacement (in bytes) into the memory region that the operation used at the target. The offset can be determined by the operation (Section 3.15) for a remote managed offset in a match list entry or by the match list entry (Section 3.12) at the target for a locally managed offset.</p> <p>At the target, this is the offset requested by the initiator.</p>
<i>uid</i>	The usage identifier of the <i>initiator</i> .
<i>initiator</i>	The identifier of the <i>initiator</i> .
<i>type</i>	Indicates the type of the full event.
<i>ptl_list</i>	The list entry or match list entry list in which the operation was delivered (See Sections 3.11.2 and 3.12.2).
<i>pt_index</i>	The portal table index where the message arrived.

ni_fail_type

Used to convey the failure of an operation. Success is indicated by PTL_NI_OK; see section 3.13.3.

atomic_operation

If this full event corresponds to an atomic operation, this indicates the atomic operation that was performed.

atomic_type

If this full event corresponds to an atomic operation, this indicates the data type of the atomic operation that was performed.

Discussion: *Notably, the full event structure does not contain a handle to the ME, LE, or MD that was associated with the full event. The `user_ptr` field is provided as the mechanism for the user to determine which ME, LE, or MD an event might be associated with.*

Table 3-3. Event Field Definition: Specification of which fields in a `ptl_event_t` structure are defined for a given event type. Fields marked with a ● are defined for both success and error conditions. Fields marked with a ○ are defined only for success conditions.

Event Type	type	initiator	pt_index	ptl_list	uid	match_bits	rlength	mlength	remote_offset	start	user_ptr	hdr_data	ni_fail_type	atomic_operation	atomic_type
PTL_EVENT_GET	●	●	●		●	●	○	●	○	●	●		●		
PTL_EVENT_GET_OVERFLOW	●	●	●		●	●	○	●	○	●	●		●		
PTL_EVENT_PUT	●	●	●		●	●	○	●	○	●	●	●	●		
PTL_EVENT_PUT_OVERFLOW	●	●	●		●	●	○	●	○	●	●	●	●		
PTL_EVENT_ATOMIC	●	●	●		●	●	○	●	○	●	●	●	●	●	●
PTL_EVENT_ATOMIC_OVERFLOW	●	●	●		●	●	○	●	○	●	●	●	●	●	●
PTL_EVENT_FETCH_ATOMIC	●	●	●		●	●	○	●	○	●	●	●	●	●	●
PTL_EVENT_FETCH_ATOMIC_OVERFLOW	●	●	●		●	●	○	●	○	●	●	●	●	●	●
PTL_EVENT_REPLY	●			○				○	○		●		●		
PTL_EVENT_SEND	●							○			●		●		
PTL_EVENT_ACK	●			○				○	○		●		●		
PTL_EVENT_PT_DISABLED	●	●											●		
PTL_EVENT_LINK	●	●									●		●		
PTL_EVENT_AUTO_UNLINK	●	●									●		●		
PTL_EVENT_AUTO_FREE	●	●									●		●		
PTL_EVENT_SEARCH	●	●	●		●	●	●	●	●	●	●	●	●	●	●
PTL_EVENT_ERROR	●										●		●		

3.13.5. PtlEQAlloc

The `PtlEQAlloc()` function is used to build an event queue. An event queue has room for at least *count* number of full events. If the event queue overflows, older events will be overwritten by new ones in most situations. If flow control is enabled on the portal table entry (See Sections 3.7.1 and 2.8) for an incoming operation, events associated with that operation will not cause an overflow, but will instead trigger a flow control event.

Function Prototype for PtIEQAlloc

```
int PtIEQAlloc(ptl_handle_ni_t ni_handle,  
              ptl_size_t count,  
              ptl_handle_eq_t *eq_handle);
```

Arguments

<i>ni_handle</i>	input	The interface handle with which the event queue will be associated.
<i>count</i>	input	A hint as to the number of full events to be stored in the event queue. An implementation may provide space for more than the requested number of event queue slots.
<i>eq_handle</i>	output	On successful return, this location will hold the newly created event queue handle.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_NO_SPACE	Indicates that there is insufficient memory to allocate the event queue.

IMPLEMENTATION NOTE 13:

Size of event queue and reserved space

Because flow control may be enabled on the portal table entries that this EQ is attached to, the implementation should insure that the space allocated for the EQ is large enough to hold the requested number of full events plus the number of portal table entries associated with this *ni_handle*. For each **PtIPTAlloc()** that enables flow control and uses a given EQ, one space should be reserved for a `PTL_EVENT_PT_DISABLED` full event associated with that EQ.

3.13.6. PtIEQFree

The **PtIEQFree()** function releases the resources associated with an event queue. It is up to the user to ensure that no memory descriptors or portal table entries are associated with the event queue before it is freed.

In the event that **PtIAbort()** was previously called, **PtIEQFree()** waits for the completion of all blocking (e.g., **PtIEQWait()**, **PtIEQPoll()**) functions related to this event queue before releasing any resources.

Once **PtIEQFree()** is called, the event queue handle stops being valid. Using the freed EQ handle in another thread results in undefined behavior.

Function Prototype for PtIEQFree

```
int PtIEQFree(ptl_handle_eq_t eq_handle);
```

Arguments

eq_handle **input** The event queue handle to be released.

Return Codes

PTL_OK Indicates success.
PTL_NO_INIT Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID Indicates that *eq_handle* is not a valid event queue handle.

3.13.7. PtlEQGet

The **PtlEQGet()** function is a non-blocking function that can be used to get the next event in an event queue. The event is removed from the queue. This function must be called with an *event* pointer to a valid **ptl_event_t** structure, which will hold the values associated with the next event in the event queue upon successful return.

Function Prototype for PtlEQGet

```
int PtlEQGet(ptl_handle_eq_t eq_handle,  
             ptl_event_t *event);
```

Arguments

eq_handle **input** The event queue handle.
event **output** On successful return, this location will hold the values associated with the next event in the event queue. *event* must point to a valid **ptl_event_t** structure.

Return Codes

PTL_OK Indicates success.
PTL_EQ_DROPPED Indicates success (i.e., an event is returned) and that at least one full event between this full event and the last full event obtained—using **PtlEQGet()**, **PtlEQWait()**, or **PtlEQPoll()**—from this event queue has been dropped due to limited space in the event queue.
PTL_NO_INIT Indicates that the Portals API has not been successfully initialized.
PTL_EQ_EMPTY Indicates that *eq_handle* is empty or another thread is waiting in **PtlEQWait()**.
PTL_ARG_INVALID Indicates that *eq_handle* is not a valid event queue handle.

3.13.8. PtlEQWait

The **PtlEQWait()** function can be used to block the calling process or thread until there is a full event in an event queue. This function returns the next event in the event queue, removes it from the queue, and returns it in the **ptl_event_t** structure passed in via the *event* pointer. In the event that multiple threads are waiting on the same event queue, **PtlEQWait()** is guaranteed to wake exactly one thread, but the order in which they are awakened is not specified. In the case that **PtlAbort()** is called, **PtlEQWait()** aborts its execution and returns.

Function Prototype for PtlEQWait

```
int PtlEQWait(ptl_handle_eq_t eq_handle,
              ptl_event_t *event);
```

Arguments

<i>eq_handle</i>	input	The event queue handle to wait on. The calling process (thread) will be blocked until the event queue is not empty.
<i>event</i>	output	On successful return, this location will hold the values associated with the next event in the event queue. <i>event</i> must point to a valid ptl_event_t structure.

Return Codes

PTL_OK	Indicates success.
PTL_EQ_DROPPED	Indicates success (i.e., an event is returned) and that at least one full event between this full event and the last full event obtained—using PtlEQGet() , PtlEQWait() , or PtlEQPoll() —from this event queue has been dropped due to limited space in the event queue.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that <i>eq_handle</i> is not a valid event queue handle.
PTL_ABORTED	Indicates that PtlAbort() was called.

3.13.9. PtlEQPoll

The **PtlEQPoll()** function can be used by the calling process to look for a full event from a set of event queues. Should an event arrive on any of the queues contained in the array of event queue handles, the full event will be removed from the given EQ and returned in the **ptl_event_t** struct passed in by *event*, which will also contain the index of the event queue from which the event was taken. In the event that multiple threads are polling the same event queue, **PtlEQPoll()** is guaranteed to wake exactly one thread, but the order in which they are awakened is not specified.

If **PtlEQPoll()** returns success, the corresponding full event is consumed. **PtlEQPoll()** provides a timeout to allow applications to poll, block for a fixed period, or block indefinitely. **PtlEQPoll()** is sufficiently general to implement both **PtlEQGet()** and **PtlEQWait()**, but these functions may allow significant optimization. If events are available on multiple event queues when **PtlEQPoll()** is called, the event returned must be from the first event queue in *eq_handles* with an available event. If all event queues are empty when **PtlEQPoll()** is called, the *timeout* is non-zero and events are inserted into multiple event queues, it is unspecified which event will be returned. **PtlEQPoll()** will also stop its execution and return when **PtlAbort()** is called.

Function Prototype for PtlEQPoll

```
int PtlEQPoll(const ptl_handle_eq_t *eq_handles,
              unsigned int size,
              ptl_time_t timeout,
              ptl_event_t *event,
              unsigned int *which);
```

Arguments

<i>eq_handles</i>	input	An array of event queue handles. All the handles must refer to the same interface although they may be from different logical network interfaces.
<i>size</i>	input	Length of the array.
<i>timeout</i>	input	Time in milliseconds to wait for a full event to occur on one of the event queue handles. The constant <code>PTL_TIME_FOREVER</code> can be used to indicate an infinite timeout.
<i>event</i>	output	On successful return (PTL_OK or PTL_EQ_DROPPED), this location will hold the values associated with the next event in the event queue. <i>event</i> must point to a valid ptl_event_t structure.
<i>which</i>	output	On successful return, this location will contain the index into <i>eq_handles</i> of the event queue from which the event was taken.

Return Codes

PTL_OK	Indicates success.
PTL_EQ_DROPPED	Indicates success (i.e., an event is returned) and that at least one full event between this full event and the last full event obtained from the event queue indicated by <i>which</i> has been dropped due to limited space in the event queue.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.
PTL_EQ_EMPTY	Indicates that the timeout has been reached and all of the event queues are empty.
PTL_ABORTED	Indicates that PtlAbort() was called.

3.14. Lightweight Counting Events

Full events copy a significant amount of data from the implementation to the application. While this data is critical for many uses (e.g. MPI), other programming models (e.g. PGAS) require very little information about individual operations. To support lightweight operations, Portals provide a lightweight event mechanism known as counting events.

Counting events are similar in semantics and event occurrence to full events (See Section 3.13.2). A counting event may be independently enabled/disabled with options on the memory descriptor, list entry, or match list entry, similar to full events. Unlike full events, counting events are disabled by default and must be explicitly enabled for a given event type. Counting events are enabled by attaching a **ptl_handle_ct_t** to a memory descriptor or list entry and

specifying which operations are to be counted in the options field. By default, counting events count the total number of operations; however, a counting event may also count the number of bytes successfully manipulated for counted operations by setting an option on the associated memory descriptor or list entry.

Counting events introduce two additional types: a user-visible representation of the counting event itself, of type `ptl_ct_event_t`, and a handle to a counting event, of type `ptl_handle_ct_t`. A counting event is allocated through a call to `PtICTAlloc()`, queried with `PtICTGet()`, `PtICTWait()`, or `PtICTPoll()`, set with `PtICTSet()`, incremented with `PtICTInc()`, and freed through a call to `PtICTFree()`. To mirror the failure semantics of the full events, counting events count success and failure events independently.

Portals event counters use unsigned integer types throughout the Portals API. Portals event counters follow the 2ⁿ modulo behavior for integer arithmetic, as specified by the ISO/IEC 10967-1 standard for language-independent arithmetic

IMPLEMENTATION

NOTE 14:

Minimizing cost of counting events

A quality implementation will attempt to minimize the cost of counting events. In many implementations, this can be done by making the `ptl_handle_ct_t` type a pointer to a `ptl_ct_event_t` structure and providing `PtICTGet()`, `PtICTWait()`, `PtICTSet()`, and `PtICTInc()` as macros which manipulate the internal structure. This may not be possible in hardware offload implementations, but `PtICTGet()` should remain as close to a pair of loads in performance as possible.

Counting events are a critical component of triggered operations, described in Section 3.16.

3.14.1. The Counting Event Type

A *ct_handle* refers to a `ptl_ct_event_t` structure. The user visible portion of this structure contains both a count of succeeding events and a count of failing events.

```
typedef struct {
    ptl_size_t success;
    ptl_size_t failure;
} ptl_ct_event_t;
```

Members

success

A count associated with successful events that counts events or bytes.

failure

A count of the number of failed events associated with the counting event.

3.14.2. PtICTAlloc

The `PtICTAlloc()` function is used to allocate a counting event that counts either operations or bytes manipulated for operations on associated memory descriptors, list entries, and match list entries. While a `PtICTAlloc()` call could be as simple as a malloc of a structure holding the counting event, it may be necessary to allocate the counting event in low memory or some other protected space. Also, it may be desirable to place all counting events in a pre-allocated array and make the *ct_handle* a simple index. A newly allocated counting event will have both the *success* and *failure* counts initialized to zero.

Function Prototype for PtlCTAlloc

```
int PtlCTAlloc(ptl_handle_ni_t ni_handle,  
              ptl_handle_ct_t *ct_handle);
```

Arguments

ni_handle **input** The interface handle with which the counting event will be associated.

ct_handle **output** On successful return, this location will hold the newly created counting event handle.

Return Codes

PTL_OK Indicates success.

PTL_NO_INIT Indicates that the Portals API has not been successfully initialized.

PTL_ARG_INVALID Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

PTL_NO_SPACE Indicates that there is insufficient memory to allocate the counting event.

3.14.3. PtlCTFree

The **PtlCTFree()** function releases the resources associated with a counting event. It is up to the user to ensure that no memory descriptors or match list entries are associated with the counting event before it is freed. On a successful return, the counting event has been released and is ready to be reallocated. As a side-effect of **PtlCTFree()**, any triggered operations waiting on the freed counting event whose thresholds have not been met will be deleted.

In the event that **PtlAbort()** was previously called, **PtlCTFree()** waits for the completion of all blocking (e.g., **PtlCTWait()**, **PtlCTPoll()**) functions related to this event queue before releasing any resources.

Once **PtlCTFree()** is called, the counting event handle stops being valid. Using the freed CT handle in another thread results in undefined behavior.

Function Prototype for PtlCTFree

```
int PtlCTFree(ptl_handle_ct_t ct_handle);
```

Arguments

ct_handle **input** The counting event handle to be released.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that <i>ct_handle</i> is not a valid counting event handle.

3.14.4. PtlCTCancelTriggered

In certain circumstances, it may be necessary to cancel triggered operations that are pending. For example, an error condition may mean that a counting event will never reach the designated threshold. **PtlCTCancelTriggered()** is provided to handle these circumstances. Upon return from **PtlCTCancelTriggered()**, all triggered operations waiting on *ct_handle* are permanently deleted. The operations are not triggered and will not modify any application-visible state. All other state associated with *ct_handle* is left unchanged.

Function Prototype for PtlCTCancelTriggered

```
int PtlCTCancelTriggered(ptl_handle_ct_t ct_handle);
```

Arguments

ct_handle **input** The counting event handle associated with the triggered operations to be canceled.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that <i>ct_handle</i> is not a valid counting event handle.

3.14.5. PtlCTGet

The **PtlCTGet()** function is used to obtain the current value of a counting event. Calling **PtlCTSet()** or **PtlCTFree()** in a separate thread while **PtlCTGet()** is executing may yield undefined results in the returned value.

Function Prototype for PtlCTGet

```
int PtlCTGet(ptl_handle_ct_t ct_handle,  
             ptl_ct_event_t *event);
```


Arguments

<i>ct_handle</i>	input	The counting event handle.
<i>event</i>	output	On successful return, this location will hold the current value associated with the counting event. <i>event</i> must point to a valid <code>ptl_ct_event_t</code> structure.

Return Codes

<code>PTL_OK</code>	Indicates success.
<code>PTL_NO_INIT</code>	Indicates that the Portals API has not been successfully initialized.
<code>PTL_ARG_INVALID</code>	Indicates that <i>ct_handle</i> is not a valid counting event handle.

3.14.6. PtlCTWait

The `PtlCTWait()` function provides blocking semantics to wait for a counting event to reach a given value. `PtlCTWait()` returns when either the *success* field of a counting event is greater than or equal to the *test* value or when the *failure* field is non-zero or when `PtlAbort()` is called. All threads that are waiting on a single counting event with a given *test* value will return from `PtlCTWait()` when that *test* value is reached.

Function Prototype for PtlCTWait

```
int PtlCTWait(ptl_handle_ct_t ct_handle,
              ptl_size_t test,
              ptl_ct_event_t *event);
```

Arguments

<i>ct_handle</i>	input	The counting event handle.
<i>test</i>	input	On successful return, the <i>success</i> field of the counting event will be greater than or equal to this value or the <i>failure</i> field of the counting event will be non-zero.
<i>event</i>	output	On successful return, this location will hold the current value associated with the counting event. <i>event</i> must point to a valid <code>ptl_ct_event_t</code> structure.

Return Codes

<code>PTL_OK</code>	Indicates success.
<code>PTL_NO_INIT</code>	Indicates that the Portals API has not been successfully initialized.
<code>PTL_ARG_INVALID</code>	Indicates that <i>ct_handle</i> is not a valid counting event handle.
<code>PTL_ABORTED</code>	Indicates that <code>PtlAbort()</code> was called.

3.14.7. PtlCTPoll

The **PtlCTPoll()** function can be used to look for one of an array of counting events where the success field has reached its respective test value. Should a counting event reach the test value for any of the counting events contained in the array of counting event handles, the value of the counting event will be returned in *event* and *which* will contain the index of the counting event from which the value was returned. **PtlCTPoll()** may be called with counting event handles associated with different logical network interfaces if those logical network interfaces all belong to the same physical network interface. **PtlCTPoll()** will also return whenever the failure field of any of the counting events is non-zero or when **PtlAbort()** is called.

PtlCTPoll() provides a timeout to allow applications to poll, block for a fixed period, or block indefinitely. If multiple counting events have reached their threshold when **PtlCTPoll()** is called, the counting event returned must be from the first counting event in *ct_handles* which has reached the threshold. If all counting events have not reached their thresholds when **PtlCTPoll()** is called, the *timeout* is non-zero and multiple counting events reach their thresholds, it is unspecified which counting event will be returned.

Function Prototype for PtlCTPoll

```
int PtlCTPoll(const ptl_handle_ct_t *ct_handles,
              const ptl_size_t *tests,
              unsigned int size,
              ptl_time_t timeout,
              ptl_ct_event_t *event,
              unsigned int *which);
```

Arguments

<i>ct_handles</i>	input	An array of counting event handles. All of the handles must refer to the same interface, although they may be from different logical network interfaces.
<i>tests</i>	input	An array of success values. PtlCTPoll() returns when any counting event in <i>ct_handles</i> would return from PtlCTWait() with the corresponding <i>test</i> in <i>tests</i> .
<i>size</i>	input	Length of the <i>ct_handles</i> and <i>tests</i> arrays.
<i>timeout</i>	input	Time in milliseconds to wait for an event to occur on one of the counting event handles. The constant PTL_TIME_FOREVER can be used to indicate an infinite timeout.
<i>event</i>	output	On successful return, this location will hold the current value associated with the counting event that caused PtlCTPoll() to return. <i>event</i> must point to a valid ptl_ct_event_t structure.
<i>which</i>	output	On successful return, this location will contain the index into <i>ct_handles</i> of the counting event that reached its test value.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates an invalid argument (e.g. a bad <i>ct_handle</i>).
PTL_CT_NONE_REACHED	Indicates that none of the counting events reached their test before the timeout was reached.
PTL_ABORTED	Indicates that PtlAbort() was called.

3.14.8. PtlCTSet

The **PtlCTSet()** function is used to set a new value for a counting event. Each field in the counting event is updated atomically relative to other updates of that field. However, there is no guarantee that the two fields are updated atomically relative to each other. The counting event must be updated before returning from **PtlCTSet()**, however the update may not be immediately visible to **PtlCTGet()**, particularly in hardware offload implementations. Both the atomicity of field updates and the delay in updating the user-visible portions of the counting event may be visible to the user, but should not affect correctness in common usage scenarios.

Function Prototype for PtlCTSet

```
int PtlCTSet(ptl_handle_ct_t ct_handle,
             ptl_ct_event_t new_ct);
```

Arguments

<i>ct_handle</i>	input	The counting event handle.
<i>new_ct</i>	input	On successful return, the value of the counting event will have been set to this value.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that <i>ct_handle</i> is not a valid counting event handle.

3.14.9. PtlCTInc

PtlCTInc() provides the ability to increment the *success* or *failure* field of a counting event. The update is atomic relative to other modifications of the counting event. To simplify implementation, the *increment* field can only be non-zero for either the *success* or *failure* field in a given call to **PtlCTInc()**. The counting event must be updated before returning from **PtlCTInc()**, however the update may not be immediately visible to **PtlCTGet()**, particularly in hardware offload implementations. This may be visible to the user, but should not affect correctness in common usage scenarios.

Function Prototype for PtlCTInc

```
int PtlCTInc(ptl_handle_ct_t ct_handle,
             ptl_ct_event_t increment);
```

Discussion: While it is possible to use negative values in calls to PtlCT* functions they will be interpreted as very large positive numbers by the Portals implementation. As such, users are discouraged from using negative values as inputs to any Portal counting event calls.

Arguments

<i>ct_handle</i>	input	The counting event handle.
<i>increment</i>	input	On successful return, the value of the counting event will have been incremented by this value.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that <i>ct_handle</i> is not a valid counting event handle.

3.15. Data Movement Operations

The Portals API provides five data movement operations: **PtlPut()**, **PtlGet()**, **PtlAtomic()**, **PtlFetchAtomic()**, and **PtlSwap()**.

3.15.1. Portals Acknowledgment Type Definition

Portals *put* and *atomic* operations which do not return data may optionally request an acknowledgment upon message delivery. Values of the type **ptl_ack_req_t** are used to specify the type of acknowledgment requested by the *initiator*. Acknowledgments are sent by the *target* when the operation has completed (i.e., when the data has been written to a list entry of the *target* process). If counting of acknowledgment events is enabled, the **PTL_MD_EVENT_CT_BYTES** option is set, and the operation is successful, the manipulated length (*mlength*) from the target is counted. If the event would indicate “failure” or the **PTL_MD_EVENT_CT_BYTES** option is not set, the number of acknowledgments is counted.

Ack Request Constants (**ptl_ack_req_t**)

PTL_ACK_REQ	An acknowledgment capable of generating both a full event and counting event is requested.
PTL_CT_ACK_REQ	An acknowledgment capable of generating a counting event is requested. Full events may or may not be generated, dependent on if an event queue is associated with the memory descriptor. If a full event is generated it is not guaranteed to contain valid data in all of the fields defined for a PTL_EVENT_ACK in table 3-3.
PTL_OC_ACK_REQ	An acknowledgment capable of generating a counting event upon operation completion is requested. An operation is considered completed when it has successfully completed Portals operation processing at the <i>target</i> . PTL_OC_ACK_REQ does not support the PTL_MD_EVENT_CT_BYTES option. The operation completion acknowledgment will indicate success as long as operation processing completed successfully. A message being dropped due to a failure to match or a permissions violation does not represent an operational failure.

PTL_NO_ACK_REQ

No acknowledgment is requested.

Discussion: *The PTL_CT_ACK_REQ and PTL_OC_ACK_REQ acknowledgment types provide significantly weaker semantics than PTL_ACK_REQ, in that the acknowledgment from the target may only contain data necessary to generate a counting event, which may improve efficiency.*

The PTL_OC_ACK_REQ acknowledgment type is useful when only operation counting is required and it is known that there is a list entry at the target that will accept the message. The PTL_OC_ACK_REQ acknowledgment type may be more efficient in some implementations because the PTL_OC_ACK_REQ acknowledgment type communicates no information about the state of the target when the message arrived. Therefore, PTL_OC_ACK_REQ may, in some implementations, be possible to implement based on transport level ACKs.

3.15.2. PtlPut

The **PtlPut()** function initiates an asynchronous *put* operation. There are several events associated with a *put* operation: completion of the send on the *initiator* node (PTL_EVENT_SEND) and the receipt of an acknowledgment (PTL_EVENT_ACK) indicating that the operation was processed by the *target*. The event PTL_EVENT_PUT is used at the *target* node to indicate the end of data delivery. In addition, PTL_EVENT_PUT_OVERFLOW can be used on the *target* node when a new entry being appended to a priority list matches a message that arrived before the corresponding match list entry had been associated with the target portal table entry (Figure 3-1 on page 76).

These (local) events will be logged using full events in the event queue or counting events in the *ct_handle* associated with the memory descriptor (*md_handle*) used in the *put* operation. Using a memory descriptor that does not have either an associated event queue or counting event results in these events being discarded. In this case, the caller must have another mechanism (e.g., a higher level protocol) for determining when it is safe to modify the memory region associated with the memory descriptor.

The local (*initiator*) offset is used to determine the starting address of the memory region within the region specified by the memory descriptor and the length specifies the length of the region in bytes. It is an error for the local offset and length parameters to specify memory outside the memory described by the memory descriptor.

Function Prototype for PtlPut

```
int PtlPut(ptl_handle_md_t md_handle,
          ptl_size_t local_offset,
          ptl_size_t length,
          ptl_ack_req_t ack_req,
          ptl_process_t target_id,
          ptl_pt_index_t pt_index,
          ptl_match_bits_t match_bits,
          ptl_size_t remote_offset,
          void *user_ptr,
          ptl_hdr_data_t hdr_data);
```

Arguments

<i>md_handle</i>	input	The memory descriptor handle that describes the memory to be sent. If the memory descriptor has an event queue associated with it, it will be used to record events when the message has been sent (PTL_EVENT_SEND, PTL_EVENT_ACK). If the memory descriptor has a counting event associated with it, it may optionally be used to record the same events.
<i>local_offset</i>	input	Offset from the start of the memory descriptor.
<i>length</i>	input	Length of the memory region to be sent.
<i>ack_req</i>	input	Controls whether an acknowledgment event is requested. Acknowledgments are only sent when they are requested by the initiating process and the memory descriptor has an event queue or counting event and the target memory descriptor enables them. If the memory descriptor is using the PTL_MD_UNRELIABLE option, then only PTL_NO_ACK_REQ may be used.
<i>target_id</i>	input	A process identifier for the <i>target</i> process.
<i>pt_index</i>	input	The index in the <i>target</i> portal table.
<i>match_bits</i>	input	The match bits to use for message selection at the <i>target</i> process (only used when matching is enabled on the network interface).
<i>remote_offset</i>	input	The offset into the target memory region (used unless the <i>target</i> match list entry has the PTL_ME_MANAGE_LOCAL option set).
<i>user_ptr</i>	input	A user-specified value that is associated with each command that can generate an event. The value does not need to be a pointer, but must fit in the space used by a pointer. This value (along with other values) is only recorded in <i>initiator</i> full events associated with this <i>put</i> operation.
<i>hdr_data</i>	input	64 bits of user data that can be included in the message header. This data is written to the full event generated at the <i>target</i> by this operation.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

Discussion: *Tying commands to a user-defined value is useful for quickly locating a user data structure associated with the *put* operation. For example, an MPI implementation can set the *user_ptr* argument to the value of an MPI Request. This direct association allows for processing of a *put* operation completion full event by the MPI implementation without a table look up or a search for the appropriate MPI Request.*

3.15.3. PtlGet

The **PtlGet()** function initiates a remote read operation. There are two events associated with a get operation. When the data is sent from the *target* node, a PTL_EVENT_GET event is registered on the *target* node if the message matched in the priority list. The message can also match in the overflow list, which will cause a PTL_EVENT_GET event to be registered on the *target* node and will later cause a PTL_EVENT_GET_OVERFLOW to be registered on the *target* node when a matching entry is appended. In either case, when the data is returned from the *target* node, a PTL_EVENT_REPLY event is registered on the *initiator* node. (Figure 3-1)

The local (*initiator*) offset is used to determine the starting address of the memory region and the length specifies the length of the region in bytes. It is an error for the local offset and length parameters to specify memory outside the memory described by the memory descriptor.

Function Prototype for PtlGet

```
int PtlGet(ptl_handle_md_t md_handle,
           ptl_size_t local_offset,
           ptl_size_t length,
           ptl_process_t target_id,
           ptl_pt_index_t pt_index,
           ptl_match_bits_t match_bits,
           ptl_size_t remote_offset,
           void *user_ptr);
```

Arguments

<i>md_handle</i>	input	The memory descriptor handle that describes the memory into which the requested data will be received. The memory descriptor can have an event queue associated with it to record full events, such as when the message receive has started. If the memory descriptor has a counting event associated with it, it may optionally be used to record the same events.
<i>local_offset</i>	input	Offset from the start of the memory descriptor.
<i>length</i>	input	Length of the memory region for the <i>reply</i> .
<i>target_id</i>	input	A process identifier for the <i>target</i> process.
<i>pt_index</i>	input	The index in the <i>target</i> portal table.
<i>match_bits</i>	input	The match bits to use for message selection at the <i>target</i> process (only used when matching is enabled on the network interface).
<i>remote_offset</i>	input	The offset into the target match list entry (used unless the target match list entry has the PTL_ME_MANAGE_LOCAL option set).
<i>user_ptr</i>	input	See the discussion for PtlPut() .

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.15.4. Portals Atomics Overview

Portals defines three closely related types of atomic operations. The **PtlAtomic()** function is a one-way operation that performs an atomic operation on data at the *target* with the data in the specified memory descriptor. The **PtlFetchAtomic()** function extends **PtlAtomic()** to be an atomic fetch-and-update operation. The value at the *target* before the operation is delivered in a *reply* message and placed into the *get* memory descriptor of the *initiator*.

Finally, the **PtlSwap()** operation atomically swaps data (including compare-and-swap and swap under mask, which require an *operand* argument).

The length of the operations performed by a **PtlAtomic()** is restricted to no more than *max_atomic_size* bytes. The *max_atomic_size* limit also guarantees that any byte in an operation (whether an atomic operation or not) that is smaller than *max_atomic_size* will only be written once in the host memory. **PtlFetchAtomic()** and **PtlSwap()** operations can be up to *max_fetch_atomic_size* bytes, except for PTL_CSWAP and PTL_MSWAP operations and their variants, which are further restricted to the length of the longest native data type.

While the length of an atomic operation is potentially multiple data items, the granularity of the atomic access is limited to the basic datatype. That is, atomic operations from different sources may be interleaved at the level of the datatype being accessed. Furthermore, atomic operations are only atomic with respect to other calls to the Portals API on the same network interface (*ni_handle*). If a network interface returned PTL_COHERENT_ATOMICS in the *features* field of **PtlInit()**, atomic operations are atomic relative to processor-initiated atomic operations, as well as any other network interface that also returned PTL_COHERENT_ATOMICS. If a network interface does not return PTL_COHERENT_ATOMICS, Portals atomic operations are not atomic relative to other host operations except those requested through the Portals API. Network interfaces that support PTL_COHERENT_ATOMICS do not need to use **PtlAtomicSync()** (a no-op for coherent interfaces), while interfaces that do not offer coherent atomics require **PtlAtomicSync()** to guarantee host visibility of data. In addition, an implementation is only required to support Portals atomic operations that are natively aligned to the size of the datatype, but it may choose to provide support for unaligned accesses. If the list entry sets the PTL_IOVEC option, a single datatype may not span multiple iovec entries. Atomicity is only guaranteed for two atomic operations using the same datatype, and overlapping atomic operations that use different datatypes are not atomic with respect to each other. The routine **PtlAtomicSync()** is provided to enable the host (or atomic operations using other datatypes) to modify memory locations that have been previously touched by an atomic operation.

The *target* match list entry must be configured to respond to *put* operations and to *get* operations if a reply is desired. The *length* argument at the initiator is used to specify the size of the request. If the *length* argument is not an integral multiple of the datatype size, the implementation may truncate to zero, truncate to a multiple of the datatype or return PTL_ARG_INVALID at the time of the call. The *mlength* in the acknowledgment will reflect the size of the actual operation performed. If the PTL_ME_MAY_ALIGN option is set, the increment of the locally managed offset may be larger than *mlength*.

There are several events that can be associated with atomic operations. When data is sent from the *initiator* node, a PTL_EVENT_SEND event is registered on the *initiator* node. It can be tracked in the event queue and/or in the counting event specified in the *put_md_handle*. The event PTL_EVENT_ATOMIC is registered on the *target* node to indicate completion of an atomic operation; and if data is returned from the *target* node, a PTL_EVENT_REPLY event is registered on the *initiator* node in the event queue and/or the counting event specified by the *get_md_handle*. Similarly, a PTL_EVENT_ACK can be registered on the *initiator* node in the event queue and/or counting event specified by the *put_md_handle* for the atomic operations that do not return data. Note that the target match list entry must have the PTL_ME_OP_PUT flag set and must also set the PTL_ME_OP_GET flag to enable a reply. As with other Portals operations, the delivery of an event indicates that the data for the associated atomic operation has been updated in application memory. This does not alleviate the requirement that all modifications of a memory location that is accessed by atomic operations must go through the Portals API. A PTL_EVENT_ATOMIC_OVERFLOW event may occur if a atomic operation matched an overflow list entry. When atomic operations match in the overflow list, the atomic operation itself is not performed on the buffer. It is the responsibility of software on the host to perform the atomic operation on the correct application memory. For atomic fetch or swap operations that match on the overflow list, the fetch operation functions in the same way as a **PtlGet()** operation that matched in the overflow list (see 3.15.3).

The three atomic functions share two new arguments introduced in Portals 4: an operation (**ptl_op_t**) and a datatype (**ptl_datatype_t**), as described below.

Discussion: To allow upper level libraries with both system defined datatype widths and fixed width datatypes to easily map to Portals, Portals provides fixed width integer types. The one exception is the long double floating-point types (PTL_LONG_DOUBLE). Because of the variability in long double encodings across systems and the lack of standard syntax for fixed width floating-point types, Portals uses a system defined width for PTL_LONG_DOUBLE and PTL_LONG_DOUBLE_COMPLEX.

Discussion: *In the case of composed atomic operations like the PTL_DIFF operation, that can be composed using a combination of two operations, a negate and sum, the event generated for matches on the priority list may not contain the PTL_DIFF operation in its optype field. For atomics that match in the unexpected list, the operation will always be listed correctly in the corresponding event. The optype field is typically used to allow for the atomic to be properly applied when the matching request is posted, for cases where the operation matched and was performed by the NIC, there is no need for the optype field in the event. When atomics can be composed from multiple different fundamental operations, the last operation performed may be inserted into the event by hardware for atomics that have successfully matched in the priority list and have been completed.*

Atomic Operation Constants (ptl_op_t)

PTL_MIN	Compute the minimum of the initiator and target value, replacing the target value with the result.
PTL_MAX	Compute the maximum of the initiator and target value, replacing the target value with the result.
PTL_SUM	Compute the sum of the initiator and target value, replacing the target value with the result.
PTL_DIFF	Compute the difference of the initiator and target value, replacing the target value with the result.
PTL_PROD	Compute the product of the initiator and target value, replacing the target value with the result.
PTL_LOR	Compute the logical OR of the initiator and target value, replacing the target value with the result.
PTL_LAND	Compute the logical AND of the initiator and target value, replacing the target value with the result.
PTL_BOR	Compute the bitwise OR of the initiator and target value, replacing the target value with the result.
PTL_BAND	Compute the bitwise AND of the initiator and target value, replacing the target value with the result.
PTL_LXOR	Compute the logical XOR of the initiator and target value, replacing the target value with the result.
PTL_BXOR	Compute the bitwise XOR of the initiator and target value, replacing the target value with the result.
PTL_SWAP	Swap the initiator and target value and deliver the target value to the buffer specified by the initiator.
PTL_CSWAP	A conditional swap. If the operand value is equal to the target value, the initiator value replaces the target value. The target value prior to comparison is always delivered to the buffer specified by the initiator. This operation is limited to single data items.
PTL_CSWAP_NE	A conditional swap. If the operand value is not equal to the target value, the initiator value replaces the target value. The target value prior to comparison is always delivered to the buffer specified by the initiator. This operation is limited to single data items.
PTL_CSWAP_LE	A conditional swap. If the operand value is less than or equal to the target value, the initiator value replaces the target value. The target value prior to comparison is always delivered to the buffer specified by the initiator. This operation is limited to single data items.

PTL_CSWAP_LT	A conditional swap. If the operand value is less than the target value, the initiator value replaces the target value. The target value prior to comparison is always delivered to the buffer specified by the initiator. This operation is limited to single data items.
PTL_CSWAP_GE	A conditional swap. If the operand value is greater than or equal to the target value, the initiator value replaces the target value. The target value prior to comparison is always delivered to the buffer specified by the initiator. This operation is limited to single data items.
PTL_CSWAP_GT	A conditional swap. If the operand value is greater than the target value, the initiator value replaces the target value. The target value prior to comparison is always delivered to the buffer specified by the initiator. This operation is limited to single data items.
PTL_MSWAP	A masked version of the swap operation. Update the bits of the target value that are set to 1 in the operand using the bits in the initiator value. Deliver the target value prior to the update to the buffer specified by the initiator. This operation is limited to single data items.

Atomic Datatype Constants (`ptl_datatype_t`)

PTL_INT8_T	8-bit signed integer
PTL_UINT8_T	8-bit unsigned integer
PTL_INT16_T	16-bit signed integer
PTL_UINT16_T	16-bit unsigned integer
PTL_INT32_T	32-bit signed integer
PTL_UINT32_T	32-bit unsigned integer
PTL_INT64_T	64-bit signed integer
PTL_UINT64_T	64-bit unsigned integer
PTL_FLOAT	32-bit floating-point number
PTL_FLOAT_COMPLEX	32-bit floating-point complex number
PTL_DOUBLE	64-bit floating-point number
PTL_DOUBLE_COMPLEX	64-bit floating-point complex number
PTL_LONG_DOUBLE	System defined long double type
PTL_LONG_DOUBLE_COMPLEX	System defined long double complex type

The legal combinations of atomic operation type, datatype, and function call are shown in Table 3-4. Generally speaking, swap operations are limited to the **PtlSwap()** function and bitwise operation are limited to integral types.

3.15.5. PtlAtomic

The **PtlAtomic()** function initiates an asynchronous *atomic* operation. The events behave like the **PtlPut()** function (see Section 3.15.2), with the exception of the target side event, which is a `PTL_EVENT_ATOMIC` (and `PTL_EVENT_ATOMIC_OVERFLOW`) instead of a `PTL_EVENT_PUT`. Similarly, the arguments mirror **PtlPut()** with the addition of a `ptl_datatype_t` and `ptl_op_t` to specify the datatype and operation being performed, respectively.

Table 3-4. Legal Atomic Operation, Datatype, and Function Combinations

	Integral Types	Floating-Point Types	Complex Types	PtlAtomic()	PtlFetchAtomic()	PtlSwap()
PTL_MIN	•	•		•	•	
PTL_MAX	•	•		•	•	
PTL_SUM	•	•	•	•	•	
PTL_PROD	•	•	•	•	•	
PTL_LOR	•			•	•	
PTL_LAND	•			•	•	
PTL_BOR	•			•	•	
PTL_BAND	•			•	•	
PTL_LXOR	•			•	•	
PTL_BXOR	•			•	•	
PTL_SWAP	•	•	•			•
PTL_CSWAP	•	•	•			•
PTL_CSWAP_NE	•	•	•			•
PTL_CSWAP_LE	•	•				•
PTL_CSWAP_LT	•	•				•
PTL_CSWAP_GE	•	•				•
PTL_CSWAP_GT	•	•				•
PTL_MSWAP	•					•
PTL_DIFF	•	•	•	•	•	

Operations performed by **PtlAtomic()** are constrained to be no more than *max_atomic_size* bytes and must be aligned at the target to the size of **ptl_datatype_t** passed in the *datatype* argument.

Function Prototype for PtlAtomic

```

int PtlAtomic(ptl_handle_md_t md_handle,
    ptl_size_t local_offset,
    ptl_size_t length,
    ptl_ack_req_t ack_req,
    ptl_process_t target_id,
    ptl_pt_index_t pt_index,
    ptl_match_bits_t match_bits,
    ptl_size_t remote_offset,
    void *user_ptr,
    ptl_hdr_data_t hdr_data,
    ptl_op_t operation,
    ptl_datatype_t datatype);

```

Arguments

- md_handle* **input** The memory descriptor handle that describes the memory to be sent. If the memory descriptor has an event queue associated with it, it will be used to record events when the message has been sent (PTL_EVENT_SEND, PTL_EVENT_ACK). If the memory descriptor has a counting event associated with it, it may optionally be used to record the same events.
- local_offset* **input** Offset from the start of the memory descriptor referenced by the *md_handle* to use for transmitted data.

<i>length</i>	input	Length of the memory region to be sent and/or received. The <i>length</i> field must be less than or equal to <i>max_atomic_size</i> .
<i>ack_req</i>	input	Controls whether an acknowledgment event is requested. Acknowledgments are only sent when they are requested by the initiating process and the memory descriptor has an event queue or counting event and the target memory descriptor enables them.
<i>target_id</i>	input	A process identifier for the <i>target</i> process.
<i>pt_index</i>	input	The index in the <i>target</i> portal table.
<i>match_bits</i>	input	The match bits to use for message selection at the <i>target</i> process.
<i>remote_offset</i>	input	The offset into the target memory region (used unless the target match list entry has the PTL_ME_MANAGE_LOCAL option set).
<i>user_ptr</i>	input	See the discussion for PtlPut() .
<i>hdr_data</i>	input	See the discussion for PtlPut() .
<i>operation</i>	input	The operation to be performed using the initiator and target data.
<i>datatype</i>	input	The type of data being operated on at the initiator and target.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.15.6. PtlFetchAtomic

The **PtlFetchAtomic()** function extends the **PtlAtomic()** function to return the value from the target *prior to the operation being performed*. When data is sent from the *initiator* node, a PTL_EVENT_SEND event is registered on the *initiator* node in the event queue and/or the counting event specified by the *put_md_handle*. The event PTL_EVENT_FETCH_ATOMIC (and potentially PTL_EVENT_FETCH_ATOMIC_OVERFLOW) is registered on the *target* node to indicate completion of an atomic operation; and if data is returned from the *target* node, a PTL_EVENT_REPLY event is registered on the *initiator* node in the event queue and/or counting event specified by the *get_md_handle*. It is an error to use memory descriptors bound to different network interfaces in a single **PtlFetchAtomic()** call. The behavior that occurs when the *local_get_offset* into the *get_md_handle* overlaps with the *local_put_offset* into the *put_md_handle* is undefined. Operations performed by **PtlFetchAtomic()** are constrained to be no more than *max_fetch_atomic_size* bytes and must be aligned at the target to the size of **ptl_datatype_t** passed in the *datatype* argument.

Function Prototype for PtlFetchAtomic

```
int PtlFetchAtomic(ptl_handle_md_t get_md_handle,
                  ptl_size_t local_get_offset,
                  ptl_handle_md_t put_md_handle,
                  ptl_size_t local_put_offset,
                  ptl_size_t length,
                  ptl_process_t target_id,
                  ptl_pt_index_t pt_index,
                  ptl_match_bits_t match_bits,
                  ptl_size_t remote_offset,
                  void *user_ptr,
                  ptl_hdr_data_t hdr_data,
                  ptl_op_t operation,
                  ptl_datatype_t datatype);
```

Arguments

<i>get_md_handle</i>	input	The memory descriptor handle that describes the memory into which the result of the operation will be placed. The memory descriptor can have an event queue associated with it to record events, such as when the result of the operation has been returned. Similarly, the memory descriptor can have a counting event to record these events.
<i>local_get_offset</i>	input	Offset from the start of the memory descriptor referenced by the <i>get_md_handle</i> to use for received data.
<i>put_md_handle</i>	input	The memory descriptor handle that describes the memory to be sent. If the memory descriptor has an event queue associated with it, it will be used to record events when the message has been sent. If the memory descriptor has a counting event associated with it, it may optionally be used to record the same events.
<i>local_put_offset</i>	input	Offset from the start of the memory descriptor referenced by the <i>put_md_handle</i> to use for transmitted data.
<i>length</i>	input	Length of the memory region to be sent and/or received. The <i>length</i> field must be less than or equal to <i>max_atomic_size</i> .
<i>target_id</i>	input	A process identifier for the <i>target</i> process.
<i>pt_index</i>	input	The index in the <i>target</i> portal table.
<i>match_bits</i>	input	The match bits to use for message selection at the <i>target</i> process.
<i>remote_offset</i>	input	The offset into the target memory region (used unless the target match list entry has the PTL_ME_MANAGE_LOCAL option set).
<i>user_ptr</i>	input	See the discussion for PtlPut() .
<i>hdr_data</i>	input	See the discussion for PtlPut() .
<i>operation</i>	input	The operation to be performed using the initiator and target data.
<i>datatype</i>	input	The type of data being operated on at the initiator and target.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.15.7. PtlSwap

The **PtlSwap()** function provides an extra argument (the *operand*) beyond the **PtlFetchAtomic()** function. Like the **PtlFetchAtomic()** function, it returns the value from the target *prior to the operation being performed*. **PtlSwap()** handles the PTL_SWAP, PTL_CSWAP (and variants), and PTL_MSWAP operations and is subject to the additional restriction that PTL_CSWAP (and variants) and PTL_MSWAP operations can only be as long as a single datatype item. Events are handled in the same way as they are for **PtlFetchAtomic()**, since **PtlSwap()** is a special case of a **PtlFetchAtomic()**. Like **PtlFetchAtomic()**, receiving a PTL_EVENT_REPLY inherently implies that the flow control check has passed on the target node. It is an error to use memory descriptors bound to different network interfaces in a single **PtlSwap()** call. The behavior that occurs when the *local_get_offset* into the *get_md_handle* overlaps with the *local_put_offset* into the *put_md_handle* is undefined. Operations performed by **PtlSwap()** are constrained to be no more than *max_fetch_atomic_size* bytes and must be aligned at the target to the size of **ptl_datatype_t** passed in the *datatype* argument. PTL_CSWAP and PTL_MSWAP operations are further restricted to one item, whose size is defined by the size of the datatype used.

Function Prototype for PtlSwap

```

int PtlSwap(ptl_handle_md_t get_md_handle,
             ptl_size_t local_get_offset,
             ptl_handle_md_t put_md_handle,
             ptl_size_t local_put_offset,
             ptl_size_t length,
             ptl_process_t target_id,
             ptl_pt_index_t pt_index,
             ptl_match_bits_t match_bits,
             ptl_size_t remote_offset,
             void *user_ptr,
             ptl_hdr_data_t hdr_data,
             const void *operand,
             ptl_op_t operation,
             ptl_datatype_t datatype);

```

Arguments

<i>get_md_handle</i>	input	The memory descriptor handle that describes the memory into which the result of the operation will be placed. The memory descriptor can have an event queue associated with it to record events, such as when the result of the operation has been returned. Similarly, the memory descriptor can have a counting event to record these events.
<i>local_get_offset</i>	input	Offset from the start of the memory descriptor referenced by the <i>get_md_handle</i> to use for received data.
<i>put_md_handle</i>	input	The memory descriptor handle that describes the memory to be sent. If the memory descriptor has an event queue associated with it, it will be used to record events when the message has been sent. If the memory descriptor has a counting event associated with it, it may optionally be used to record the same events.

<i>local_put_offset</i>	input	Offset from the start of the memory descriptor referenced by the <i>put_md_handle</i> to use for transmitted data.
<i>length</i>	input	Length of the memory region to be sent and/or received. The <i>length</i> field must be less than or equal to <i>max_atomic_size</i> for PTL_SWAP operations and can only be as large as a single datatype item for PTL_CSWAP and PTL_MSWAP operations, and variants of those.
<i>target_id</i>	input	A process identifier for the <i>target</i> process.
<i>pt_index</i>	input	The index in the <i>target</i> portal table.
<i>match_bits</i>	input	The match bits to use for message selection at the <i>target</i> process.
<i>remote_offset</i>	input	The offset into the target memory region (used unless the target match list entry has the PTL_ME_MANAGE_LOCAL option set).
<i>user_ptr</i>	input	See the discussion for PtlPut() .
<i>hdr_data</i>	input	See the discussion for PtlPut() .
<i>operand</i>	input	A pointer to the data to be used for the PTL_CSWAP (and variants) and PTL_MSWAP operations (ignored for other operations). The data pointed to is of the type specified by the <i>datatype</i> argument and must be included in the message.
<i>operation</i>	input	The operation to be performed using the initiator and target data.
<i>datatype</i>	input	The type of data being operated on at the initiator and target.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.15.8. PtlAtomicSync

The **PtlAtomicSync()** function ensures host visibility of the atomic accesses previously completed through the Portals API. When a data item is accessed by a Portals atomic operation, modification of the same data item by the host or by an atomic operation using a different datatype can lead to undefined behavior. When **PtlAtomicSync()** is called, it will block until it is safe for the host (or other atomic operations with a different datatype) to modify the data items touched by previous Portals atomic operations. **PtlAtomicSync()** is called at the target of atomic operations. For NIs that provide PTL_COHERENT_ATOMICS calls to **PtlAtomicSync()** are unnecessary.

IMPLEMENTATION NOTE 15: Portals Atomic Synchronization

The atomicity definition for Portals allows a network interface to offload atomic operations and to have a non-coherent cache on the network interface. With a non-coherent cache, any access to a memory location by an atomic operation makes it impossible to safely modify that location on the host. **PtlAtomicSync()** is provided to make modifications from the host safe again.

Function Prototype for PtlAtomicSync

```
int PtlAtomicSync();
```

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.

3.16. Triggered Operations

For a variety of scenarios, it is desirable to setup a response to incoming messages. As an example, a tree based reduction operation could be performed by having each layer of the tree issue a **PtlAtomic()** operation to its parent after receiving a **PtlAtomic()** from all of its children. To provide this operation, triggered versions of each of the data movement operations are provided. To create a triggered operation, a *trig_ct_handle* and an integer *threshold* are added to the argument list. When the success field of the count (not including failures) referenced by the *trig_ct_handle* argument reaches or exceeds the *threshold* (equal to or greater), the operation proceeds *at the initiator of the operation*. For example, a **PtlTriggeredGet()** or a **PtlTriggeredAtomic()** will not leave the *initiator* until the threshold is reached. A triggered operation does not use the state of the buffer when the application calls the Portals function. Instead, it uses the state of the buffer after the threshold condition is met. Pending triggered operations can be canceled using **PtlCTCancelTriggered()**.

Triggered operations are processed in order of threshold values, even if the counting event is increased by a large amount at once (such as through a call to **PtlCTInc()**). If a counting event has already reached the *threshold* when a triggered operation is created, that operation is immediately processed.

Triggered operations proceed in the order their trigger threshold is reached, implying ordering within the implementation.

Discussion: *The use of a *trig_ct_handle* and *threshold* enables a variety of usage models. A single match list entry can trigger one operation (or several) by using an independent *trig_ct_handle* on the match list entry. One operation can be triggered by a combination of previous events (include a combination of initiator and target side events) by having all of the earlier operations reference a single *trig_ct_handle* and using an appropriate threshold. Users are strongly encouraged to order calls to triggered operations by increasing threshold value as there may be significant performance advantages to ordering calls this way.*

IMPLEMENTATION NOTE 16:

Ordering of Triggered Operations

The semantics of triggered operations imply that (at a minimum) operations will proceed in the order that their trigger threshold is reached. An implementation will release operations that reach their threshold simultaneously on the same *trig_ct_handle* in the order that they are issued. Users should also create triggered operations in ascending *threshold* values to decrease sorting work on implementations.

3.16.1. PtlTriggeredPut

The **PtlTriggeredPut()** function adds triggered operation semantics to the **PtlPut()** function described in Section 3.15.2.

Function Prototype for PtlTriggeredPut

```
int PtlTriggeredPut(ptl_handle_md_t md_handle,
                   ptl_size_t local_offset,
                   ptl_size_t length,
                   ptl_ack_req_t ack_req,
                   ptl_process_t target_id,
                   ptl_pt_index_t pt_index,
                   ptl_match_bits_t match_bits,
                   ptl_size_t remote_offset,
                   void *user_ptr,
                   ptl_hdr_data_t hdr_data,
                   ptl_handle_ct_t trig_ct_handle,
                   ptl_size_t threshold);
```

Arguments

<i>md_handle</i>	input	See PtlPut() description in Section 3.15.2.
<i>local_offset</i>	input	See PtlPut() description in Section 3.15.2.
<i>length</i>	input	See PtlPut() description in Section 3.15.2.
<i>ack_req</i>	input	See PtlPut() description in Section 3.15.2.
<i>target_id</i>	input	See PtlPut() description in Section 3.15.2.
<i>pt_index</i>	input	See PtlPut() description in Section 3.15.2.
<i>match_bits</i>	input	See PtlPut() description in Section 3.15.2.
<i>remote_offset</i>	input	See PtlPut() description in Section 3.15.2.
<i>user_ptr</i>	input	See PtlPut() description in Section 3.15.2.
<i>hdr_data</i>	input	See PtlPut() description in Section 3.15.2.
<i>trig_ct_handle</i>	input	Handle used for triggering the operation.
<i>threshold</i>	input	Threshold at which the operation triggers.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_NO_SPACE	Indicates there is insufficient memory to register the triggered operation.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.16.2. PtlTriggeredGet

The **PtlTriggeredGet()** function adds triggered operation semantics to the **PtlGet()** function described in Section 3.15.3.

Function Prototype for PtlTriggeredGet

```
int PtlTriggeredGet(ptl_handle_md_t md_handle,
                   ptl_size_t local_offset,
                   ptl_size_t length,
                   ptl_process_t target_id,
                   ptl_pt_index_t pt_index,
                   ptl_match_bits_t match_bits,
                   ptl_size_t remote_offset,
                   void *user_ptr,
                   ptl_handle_ct_t ct_handle,
                   ptl_size_t threshold);
```

Arguments

<i>md_handle</i>	input	See PtlGet() description in Section 3.15.3.
<i>local_offset</i>	input	See PtlGet() description in Section 3.15.3.
<i>length</i>	input	See PtlGet() description in Section 3.15.3.
<i>target_id</i>	input	See PtlGet() description in Section 3.15.3.
<i>pt_index</i>	input	See PtlGet() description in Section 3.15.3.
<i>match_bits</i>	input	See PtlGet() description in Section 3.15.3.
<i>remote_offset</i>	input	See PtlGet() description in Section 3.15.3.
<i>user_ptr</i>	input	See PtlGet() description in Section 3.15.3.
<i>trig_ct_handle</i>	input	Handle used for triggering the operation.
<i>threshold</i>	input	Threshold at which the operation triggers.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_NO_SPACE	Indicates there is insufficient memory to register the triggered operation.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.16.3. PtlTriggeredAtomic

The **PtlTriggeredAtomic()** function adds triggered operation semantics to the **PtlAtomic()** function described in Section 3.15.5. When combined with triggered counting increments (**PtlTriggeredCTInc()**) and sets

(**PtlTriggeredCTSet()**), triggered atomic operations enable an offloaded, non-blocking implementation of most collective operations.

Function Prototype for PtlTriggeredAtomic

```
int PtlTriggeredAtomic(ptl_handle_md_t md_handle,
                       ptl_size_t local_offset,
                       ptl_size_t length,
                       ptl_ack_req_t ack_req,
                       ptl_process_t target_id,
                       ptl_pt_index_t pt_index,
                       ptl_match_bits_t match_bits,
                       ptl_size_t remote_offset,
                       void *user_ptr,
                       ptl_hdr_data_t hdr_data,
                       ptl_op_t operation,
                       ptl_datatype_t datatype,
                       ptl_handle_ct_t trig_ct_handle,
                       ptl_size_t threshold);
```

Arguments

<i>md_handle</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>local_offset</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>length</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>ack_req</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>target_id</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>pt_index</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>match_bits</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>remote_offset</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>user_ptr</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>hdr_data</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>operation</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>datatype</i>	input	See PtlAtomic() description in Section 3.15.5.
<i>trig_ct_handle</i>	input	Handle used for triggering the operation.
<i>threshold</i>	input	Threshold at which the operation triggers.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_NO_SPACE	Indicates there is insufficient memory to register the triggered operation.

PTL_ARG_INVALID

Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.16.4. PtlTriggeredFetchAtomic

The **PtlTriggeredFetchAtomic()** function adds triggered operation semantics to the **PtlFetchAtomic()** function described in Section 3.15.6.

Function Prototype for PtlTriggeredFetchAtomic

```
int PtlTriggeredFetchAtomic(ptl_handle_md_t get_md_handle,
                             ptl_size_t local_get_offset,
                             ptl_handle_md_t put_md_handle,
                             ptl_size_t local_put_offset,
                             ptl_size_t length,
                             ptl_process_t target_id,
                             ptl_pt_index_t pt_index,
                             ptl_match_bits_t match_bits,
                             ptl_size_t remote_offset,
                             void *user_ptr,
                             ptl_hdr_data_t hdr_data,
                             ptl_op_t operation,
                             ptl_datatype_t datatype,
                             ptl_handle_ct_t trig_ct_handle,
                             ptl_size_t threshold);
```

Arguments

<i>get_md_handle</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>local_get_offset</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>put_md_handle</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>local_put_offset</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>length</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>target_id</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>pt_index</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>match_bits</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>remote_offset</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>user_ptr</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>hdr_data</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>operation</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>datatype</i>	input	See PtlFetchAtomic() description in Section 3.15.6.
<i>trig_ct_handle</i>	input	Handle used for triggering the operation.
<i>threshold</i>	input	Threshold at which the operation triggers.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_NO_SPACE	Indicates there is insufficient memory to register the triggered operation.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.16.5. PtlTriggeredSwap

The **PtlTriggeredSwap()** function adds triggered operation semantics to the **PtlSwap()** function described in Section 3.15.7.

Function Prototype for PtlTriggeredSwap

```
int PtlTriggeredSwap(ptl_handle_md_t get_md_handle,  
                    ptl_size_t local_get_offset,  
                    ptl_handle_md_t put_md_handle,  
                    ptl_size_t local_put_offset,  
                    ptl_size_t length,  
                    ptl_process_t target_id,  
                    ptl_pt_index_t pt_index,  
                    ptl_match_bits_t match_bits,  
                    ptl_size_t remote_offset,  
                    void *user_ptr,  
                    ptl_hdr_data_t hdr_data,  
                    const void *operand,  
                    ptl_op_t operation,  
                    ptl_datatype_t datatype,  
                    ptl_handle_ct_t trig_ct_handle,  
                    ptl_size_t threshold);
```

Arguments

<i>get_md_handle</i>	input	See PtlSwap() description in Section 3.15.7.
<i>local_get_offset</i>	input	See PtlSwap() description in Section 3.15.7.
<i>put_md_handle</i>	input	See PtlSwap() description in Section 3.15.7.
<i>local_put_offset</i>	input	See PtlSwap() description in Section 3.15.7.
<i>length</i>	input	See PtlSwap() description in Section 3.15.7.
<i>target_id</i>	input	See PtlSwap() description in Section 3.15.7.
<i>pt_index</i>	input	See PtlSwap() description in Section 3.15.7.
<i>match_bits</i>	input	See PtlSwap() description in Section 3.15.7.
<i>remote_offset</i>	input	See PtlSwap() description in Section 3.15.7.
<i>user_ptr</i>	input	See PtlSwap() description in Section 3.15.7.

<i>hdr_data</i>	input	See PtlSwap() description in Section 3.15.7.
<i>operand</i>	input	See PtlSwap() description in Section 3.15.7.
<i>operation</i>	input	See PtlSwap() description in Section 3.15.7.
<i>datatype</i>	input	See PtlSwap() description in Section 3.15.7.
<i>trig_ct_handle</i>	input	Handle used for triggering the operation.
<i>threshold</i>	input	Threshold at which the operation triggers.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_NO_SPACE	Indicates there is insufficient memory to register the triggered operation.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.16.6. PtlTriggeredCTInc

The triggered counting event increment extends the counting event increment (**PtlCTInc()**) with the triggered operation semantics. It is a convenient mechanism to provide chaining of dependencies between counting events. This allows a relatively arbitrary ordering of operations. For example, a **PtlTriggeredPut()** and a **PtlTriggeredCTInc()** could be dependent on *ct_handle* A with the same threshold. If the **PtlTriggeredCTInc()** is set to increment *ct_handle* B and a second **PtlTriggeredPut()** is dependent on *ct_handle* B, the second **PtlTriggeredPut()** will occur after the first.

Function Prototype for PtlTriggeredCTInc

```
int PtlTriggeredCTInc(ptl_handle_ct_t ct_handle,
                    ptl_ct_event_t increment,
                    ptl_handle_ct_t trig_ct_handle,
                    ptl_size_t threshold);
```

Arguments

<i>ct_handle</i>	input	See PtlCTInc() description in Section 3.14.9.
<i>increment</i>	input	See PtlCTInc() description in Section 3.14.9.
<i>trig_ct_handle</i>	input	Handle used for triggering the operation.
<i>threshold</i>	input	Threshold at which the operation triggers.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_NO_SPACE	Indicates there is insufficient memory to register the triggered operation.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.16.7. PtlTriggeredCTSet

The triggered counting event increment extends the counting event set (**PtICTSet()**) with the triggered operation semantics. It is a convenient mechanism to provide reinitialization of counting events between invocations of an algorithm.

Function Prototype for PtlTriggeredCTSet

```
int PtlTriggeredCTSet(ptl_handle_ct_t ct_handle,
                    ptl_ct_event_t new_ct,
                    ptl_handle_ct_t trig_ct_handle,
                    ptl_size_t threshold);
```

Arguments

<i>ct_handle</i>	input	See PtICTSet() description in Section 3.14.8.
<i>new_ct</i>	input	See PtICTSet() description in Section 3.14.8.
<i>trig_ct_handle</i>	input	Handle used for triggering the operation.
<i>threshold</i>	input	Threshold at which the operation triggers.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_NO_SPACE	Indicates there is insufficient memory to register the triggered operation.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.17. Deferred Communication Operations

Frequently, upper layer protocols and applications generate a stream of operations with loose synchronization requirements between operations. For example, an MPI implementation may need to start a large number of operations to implement the fan-out portion of a collective operation. The portals deferred communication operations provide a mechanism for allowing the Portals implementation to optimize for these situations.

3.17.1. PtlStartBundle

The **PtlStartBundle()** function is used by the application to indicate to the implementation that a group of communication operations is about to start. **PtlStartBundle()** takes an *ni_handle* as an argument and only impacts operations on that *ni_handle*. **PtlStartBundle()** can be called multiple times, and each call to **PtlStartBundle()** increments a reference count and must be matched by a call to **PtlEndBundle()**. After a call to **PtlStartBundle()**, the implementation may begin deferring communication operations until a call to **PtlEndBundle()**.

Function Prototype for PtlStartBundle

```
int PtlStartBundle(ptl_handle_ni_t ni_handle);
```

Arguments

ni_handle **input** An interface handle to start bundling operations.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

Discussion: *Layered libraries and heavily nested **PtlStartBundle()** calls can yield unexpected results. The **PtlStartBundle()** and **PtlEndBundle()** interface was designed for use in short periods of high activity (e.g. during the setup of a collective operation or during an inner loop for PGAS languages). The interval between **PtlStartBundle()** and the corresponding **PtlEndBundle()** should be kept short.*

IMPLEMENTATION NOTE 17:

Purpose of Bundling

The **PtlStartBundle()** and **PtlEndBundle()** interface was designed to allow the implementation to avoid unnecessary `sence()`/memory barrier operations during periods that the application expects high message rate. A quality implementation will attempt to minimize latency while maximizing message rate. For example, an implementation that requires writes into “write-combining” space may require `sence()` operations with every message to have relatively deterministic latency. Between a **PtlStartBundle()** and **PtlEndBundle()**, the implementation might simply omit the `sence()` operations.

3.17.2. PtlEndBundle

The **PtlEndBundle()** function is used by the application to indicate to the implementation that a group of communication operations has ended. **PtlEndBundle()** takes an *ni_handle* as an argument and only impacts operations on that *ni_handle*. **PtlEndBundle()** must be called once for each **PtlStartBundle()** call. At each call to **PtlEndBundle()**, the implementation must initiate all communication operations that have been deferred; however, the implementation is not required to cease bundling future operations until the reference count reaches zero.

Function Prototype for PtlEndBundle

```
int PtlEndBundle(ptl_handle_ni_t ni_handle);
```

Arguments

ni_handle **input** An interface handle to end bundling operations.

Return Codes

PTL_OK	Indicates success.
PTL_NO_INIT	Indicates that the Portals API has not been successfully initialized.
PTL_ARG_INVALID	Indicates that an invalid argument was passed. The definition of which arguments are checked is implementation dependent.

3.18. Operations on Handles

Handles are opaque data types. The only operation defined on them by the Portals API is a comparison function.

3.18.1. PtlHandleIsEqual

The **PtlHandleIsEqual()** function compares two handles to determine if they represent the same object. **PtlHandleIsEqual()** does not check whether the two handles are valid, but only whether they are equal.

Function Prototype for PtlHandleIsEqual

```
int PtlHandleIsEqual(ptl_handle_any_t handle1,  
                    ptl_handle_any_t handle2);
```

Arguments

<i>handle1</i>	input	An object handle. May be the constant value <code>PTL_INVALID_HANDLE</code> , which represents the value of an invalid handle.
<i>handle2</i>	input	An object handle. May be the constant value <code>PTL_INVALID_HANDLE</code> , which represents the value of an invalid handle.

Return Codes

zero	Indicates that the two handles are not equivalent.
non-zero	Indicates that the two handles are equivalent.

Discussion: *PtlHandleIsEqual()* returns a value suitable for direct evaluation in a conditional expression. While different from all other Portals functions and previous Portals versions, it does greatly simplify usage of *PtlHandleIsEqual()*.

3.19. Summary

We conclude this chapter by summarizing the names introduced by the Portals API. We start with the data types introduced by the API. This is followed by a summary of the functions defined by the API which is followed by a summary of the function return codes. Finally, we conclude with a summary of the other constant values defined by the API.

Table 3-5 presents a summary of the types defined by the Portals API. The first column in this table gives the type name, the second column gives a brief description of the type, the third column identifies the section where the type is defined, and the fourth column lists the functions that have arguments of this type and structures with members of this type.

Table 3-5. Portals Data Types: Data Types Defined by the Portals API.

Name	Meaning	Definition	Functions/Data Structures
<code>ptl_ack_req_t</code>	acknowledgment request types	3.15.1	<code>PtlAtomic()</code> , <code>PtlPut()</code> , <code>PtlTriggeredAtomic()</code> , <code>PtlTriggeredPut()</code>
<code>ptl_ct_event_t</code>	counting event structure	3.14.1	<code>PtlCTGet()</code> , <code>PtlCTInc()</code> , <code>PtlCTPoll()</code> , <code>PtlCTSet()</code> , <code>PtlTriggeredCTInc()</code> , <code>PtlTriggeredCTSet()</code> , <code>PtlCTWait()</code>
<code>ptl_datatype_t</code>	datatype for atomic operation	3.15.4	<code>PtlAtomic()</code> , <code>PtlFetchAtomic()</code> , <code>PtlSwap()</code> , <code>PtlTriggeredAtomic()</code> , <code>PtlTriggeredFetchAtomic()</code> , <code>PtlTriggeredSwap()</code> , <code>ptl_event_t</code>
<code>ptl_event_kind_t</code>	event kind	3.13.1	<code>ptl_event_t</code>
<code>ptl_event_t</code>	event queue entry	3.13.4	<code>PtlEQGet()</code> , <code>PtlEQWait()</code> , <code>PtlEQPoll()</code>
<code>ptl_handle_any_t</code>	any object handles	3.3.2	<code>PtlHandleIsEqual()</code> , <code>PtlNIHandle()</code>
<code>ptl_handle_ct_t</code>	counting event handles	3.3.2	<code>PtlCTAlloc()</code> , <code>PtlCTCancelTriggered()</code> , <code>PtlCTFree()</code> , <code>PtlCTGet()</code> , <code>PtlCTInc()</code> , <code>PtlCTPoll()</code> , <code>PtlCTSet()</code> , <code>PtlCTWait()</code> , <code>PtlTriggeredAtomic()</code> , <code>PtlTriggeredCTInc()</code> , <code>PtlTriggeredCTSet()</code> , <code>PtlTriggeredFetchAtomic()</code> , <code>PtlTriggeredGet()</code> , <code>PtlTriggeredPut()</code> , <code>PtlTriggeredSwap()</code> , <code>ptl_le_t</code> , <code>ptl_md_t</code> , <code>ptl_me_t</code>

continued on next page

continued from previous page

Name	Meaning	Definition	Functions/Data Structures
<code>ptl_handle_eq_t</code>	event queue handles	3.3.2	<code>PtlEQAlloc()</code> , <code>PtlEQFree()</code> , <code>PtlEQGet()</code> , <code>PtlEQPoll()</code> , <code>PtlEQWait()</code> , <code>PtlPTAlloc()</code> , <code>ptl_md_t</code>
<code>ptl_handle_le_t</code>	list entry handles	3.3.2	<code>PtlLEAppend()</code> , <code>PtlLEUnlink()</code>
<code>ptl_handle_md_t</code>	memory descriptor handles	3.3.2	<code>PtlAtomic()</code> , <code>PtlFetchAtomic()</code> , <code>PtlGet()</code> , <code>PtlMDBind()</code> , <code>PtlMDRelease()</code> , <code>PtlPut()</code> , <code>PtlSwap()</code> , <code>PtlTriggeredAtomic()</code> , <code>PtlTriggeredFetchAtomic()</code> , <code>PtlTriggeredGet()</code> , <code>PtlTriggeredPut()</code> , <code>PtlTriggeredSwap()</code>
<code>ptl_handle_me_t</code>	match list entry handles	3.3.2	<code>PtlMEAppend()</code> , <code>PtlMEUnlink()</code>
<code>ptl_handle_ni_t</code>	network interface handles	3.3.2	<code>PtlCTAlloc()</code> , <code>PtlEQAlloc()</code> , <code>PtlEndBundle()</code> , <code>PtlGetId()</code> , <code>PtlGetMap()</code> , <code>PtlGetPhysId()</code> , <code>PtlGetUid()</code> , <code>PtlLEAppend()</code> , <code>PtlLESearch()</code> , <code>PtlMDBind()</code> , <code>PtlMEAppend()</code> , <code>PtlMESearch()</code> , <code>PtlNIFini()</code> , <code>PtlNIHandle()</code> , <code>PtlNIInit()</code> , <code>PtlNIStatus()</code> , <code>PtlPTAlloc()</code> , <code>PtlPTDisable()</code> , <code>PtlPTEnable()</code> , <code>PtlPTFree()</code> , <code>PtlSetMap()</code> , <code>PtlStartBundle()</code>
<code>ptl_hdr_data_t</code>	user header data	3.15.2	<code>PtlAtomic()</code> , <code>PtlFetchAtomic()</code> , <code>PtlPut()</code> , <code>PtlSwap()</code> , <code>PtlTriggeredAtomic()</code> , <code>PtlTriggeredFetchAtomic()</code> , <code>PtlTriggeredPut()</code> , <code>PtlTriggeredSwap()</code> , <code>ptl_event_t</code>
<code>ptl_interface_t</code>	network interface identifiers	3.3.5	<code>PtlNIInit()</code>
<code>ptl_iovec_t</code>	scatter/gather buffer descriptors	3.10.2	
<code>ptl_le_t</code>	list entries	3.11.1	<code>PtlLEAppend()</code> , <code>PtlLESearch()</code>
<code>ptl_list_t</code>	type of list attached to a portal table entry	3.12.2	<code>PtlLEAppend()</code> , <code>PtlLEAppend()</code> , <code>ptl_event_t</code>
<code>ptl_match_bits_t</code>	match (and ignore) bits	3.3.4	<code>PtlAtomic()</code> , <code>PtlFetchAtomic()</code> , <code>PtlGet()</code> , <code>PtlPut()</code> , <code>PtlSwap()</code> , <code>PtlTriggeredAtomic()</code> , <code>PtlTriggeredFetchAtomic()</code> , <code>PtlTriggeredGet()</code> , <code>PtlTriggeredPut()</code> , <code>PtlTriggeredSwap()</code> , <code>ptl_event_t</code> , <code>ptl_me_t</code>
<code>ptl_md_t</code>	memory descriptors	3.10.1	<code>PtlMDBind()</code>
<code>ptl_me_t</code>	match list entries	3.12.1	<code>PtlMEAppend()</code> , <code>PtlMESearch()</code>
<code>ptl_ni_fail_t</code>	network interface specific failures	3.13.3	<code>ptl_event_t</code>
<code>ptl_ni_limits_t</code>	implementation dependent limits	3.6.1	<code>PtlNIInit()</code>
<code>ptl_nid_t</code>	node identifiers	3.3.6	<code>ptl_process_t</code>
<code>ptl_op_t</code>	atomic operation type	3.15.4	<code>PtlAtomic()</code> , <code>PtlFetchAtomic()</code> , <code>PtlSwap()</code> , <code>PtlTriggeredAtomic()</code> , <code>PtlTriggeredFetchAtomic()</code> , <code>PtlTriggeredSwap()</code> , <code>ptl_event_t</code>
<code>ptl_pid_t</code>	process identifier	3.3.6	<code>PtlNIInit()</code> , <code>ptl_process_t</code>

continued on next page

continued from previous page

Name	Meaning	Definition	Functions/Data Structures
<code>ptl_process_t</code>	process identifiers	3.9.1	<code>PtlAtomic()</code> , <code>PtlFetchAtomic()</code> , <code>PtlGet()</code> , <code>PtlGetId()</code> , <code>PtlGetMap()</code> , <code>PtlGetPhysId()</code> , <code>PtlPut()</code> , <code>PtlSetMap()</code> , <code>PtlSwap()</code> , <code>PtlTriggeredAtomic()</code> , <code>PtlTriggeredFetchAtomic()</code> , <code>PtlTriggeredGet()</code> , <code>PtlTriggeredPut()</code> , <code>PtlTriggeredSwap()</code> , <code>ptl_event_t</code> , <code>ptl_me_t</code>
<code>ptl_pt_index_t</code>	portal table indexes	3.3.3	<code>PtlAtomic()</code> , <code>PtlFetchAtomic()</code> , <code>PtlGet()</code> , <code>PtlLEAppend()</code> , <code>PtlLESearch()</code> , <code>PtlIMEAppend()</code> , <code>PtlIMESearch()</code> , <code>PtlPTAlloc()</code> , <code>PtlPTDisable()</code> , <code>PtlPTEnable()</code> , <code>PtlPTFree()</code> , <code>PtlPut()</code> , <code>PtlSwap()</code> , <code>PtlTriggeredAtomic()</code> , <code>PtlTriggeredFetchAtomic()</code> , <code>PtlTriggeredGet()</code> , <code>PtlTriggeredPut()</code> , <code>PtlTriggeredSwap()</code> , <code>ptl_event_t</code> <code>ptl_process_t</code>
<code>ptl_rank_t</code>	rank within a group of communicating processes	3.3.6	
<code>ptl_search_op_t</code>	operation performed by list search	3.12.4	<code>PtlLESearch()</code> , <code>PtlIMESearch()</code>
<code>ptl_size_t</code>	sizes	3.3.1	<code>PtlAtomic()</code> , <code>PtlCTPoll()</code> , <code>PtlCTWait()</code> , <code>PtlEQAlloc()</code> , <code>PtlFetchAtomic()</code> , <code>PtlGet()</code> , <code>PtlGetMap()</code> , <code>PtlPut()</code> , <code>PtlSetMap()</code> , <code>PtlSwap()</code> , <code>PtlTriggeredAtomic()</code> , <code>PtlTriggeredCTInc()</code> , <code>PtlTriggeredCTSet()</code> , <code>PtlTriggeredFetchAtomic()</code> , <code>PtlTriggeredGet()</code> , <code>PtlTriggeredPut()</code> , <code>PtlTriggeredSwap()</code> , <code>ptl_ct_event_t</code> , <code>ptl_event_t</code> , <code>ptl_iovec_t</code> , <code>ptl_le_t</code> , <code>ptl_md_t</code> , <code>ptl_me_t</code> , <code>ptl_ni_limits_t</code>
<code>ptl_sr_index_t</code>	status register indexes	3.3.7	<code>PtlNIStatus()</code>
<code>ptl_sr_value_t</code>	status register values	3.3.7	<code>PtlNIStatus()</code>
<code>ptl_time_t</code>	time in milliseconds	3.13.9	<code>PtlCTPoll()</code> , <code>PtlEQPoll()</code>
<code>ptl_uid_t</code>	usage identifier	3.3.6	<code>PtlGetUid()</code> , <code>ptl_event_t</code> , <code>ptl_le_t</code> , <code>ptl_me_t</code>

Table 3-6 presents a summary of the functions defined by the Portals API. The first column in this table gives the name for the function, the second column gives a brief description of the operation implemented by the function, and the third column identifies the section where the function is defined.

Table 3-6. Portals Functions: Functions Defined by the Portals API.

Name	Meaning	Definition
<code>PtlAbort()</code>	abort polling and waiting calls	3.5.3
<code>PtlAtomic()</code>	perform an atomic operation	3.15.5
<code>PtlAtomicSync()</code>	synchronize results of atomic operations with the host	3.15.8
<code>PtlCTAlloc()</code>	create a counting event	3.14.2

continued on next page

Name	Meaning	Definition
PtICTCancelTriggered()	cancel pending triggered operations	3.14.4
PtICTFree()	free a counting event	3.14.3
PtICTGet()	get the current value of a counting event	3.14.5
PtICTInc()	increment a counting event by a certain value	3.14.9
PtICTPoll()	wait for an array of counting events to reach certain values	3.14.7
PtICTSet()	set a counting event to a certain value	3.14.8
PtICTWait()	wait for a counting event to reach a certain value	3.14.6
PtIEndBundle()	end a communications bundle	3.17.2
PtIEQAlloc()	create an event queue	3.13.5
PtIEQFree()	release the resources for an event queue	3.13.6
PtIEQGet()	get the next event from an event queue	3.13.7
PtIEQPoll()	poll for a new event on multiple event queues	3.13.9
PtIEQWait()	wait for a new event in an event queue	3.13.8
PtIFetchAtomic()	perform an fetch and atomic operation	3.15.6
PtIFini()	shut down the Portals API	3.5.2
PtIGet()	perform a <i>get</i> operation	3.15.3
PtIGetId()	get the identifier for the current process	3.9.2
PtIGetMap()	retrieve a rank to physical mapping	3.6.7
PtIGetPhysId()	get the physical identifier for the current process	3.9.3
PtIGetUid()	get the network interface specific usage identifier	3.8.1
PtIHandleIsEqual()	compares two handles to determine if they represent the same object	3.18.1
PtIInit()	initialize the Portals API	3.5.1
PtILEAppend()	create a list entry and append it to a portal table	3.11.2
PtILESearch()	search an unexpected header	3.11.4
PtILEUnlink()	remove a list entry from a list and release its resources	3.11.3
PtIMDBind()	create a free-floating memory descriptor	3.10.3
PtIMDRelease()	release resources associated with a memory descriptor	3.10.4
PtIMEAppend()	create a match list entry and append it to a portal table	3.12.2
PtIMESearch()	search an unexpected header	3.12.4
PtIMEUnlink()	remove a match list entry from a list and release its resources	3.12.3
PtINIFini()	shut down a network interface	3.6.3
PtINIHandle()	get the network interface handle for an object	3.6.5
PtINIInit()	initialize a network interface	3.6.2
PtINIStatus()	read a network interface status register	3.6.4
PtIPTAlloc()	allocate a free portal table entry	3.7.1
PtIPTFree()	free a portal table entry	3.7.2
PtIPTDisable()	disable a portal table entry	3.7.3
PtIPTEnable()	enable a portal table entry that has been disabled	3.7.4
PtIPut()	perform a <i>put</i> operation	3.15.2
PtISetMap()	initialize a rank to physical mapping	3.6.6
PtIStartBundle()	start a communications bundle	3.17.1
PtISwap()	perform a swap operation	3.15.7
PtITriggeredAtomic()	perform a triggered atomic operation	3.16.3
PtITriggeredCTInc()	a triggered increment of a counting event by a certain value	3.16.6
PtITriggeredCTSet()	a triggered set of a counting event by a certain value	3.16.7
PtITriggeredFetchAtomic()	perform a triggered fetch and atomic operation	3.16.4
PtITriggeredGet()	perform a triggered <i>get</i> operation	3.16.2
PtITriggeredPut()	perform a triggered <i>put</i> operation	3.16.1
PtITriggeredSwap()	perform a triggered swap operation	3.16.5

Table 3-7 summarizes the return codes used by functions defined by the Portals API. The first column of this table

gives the symbolic name for the constant, the second column gives a brief description of the value, and the third column identifies the functions that can return this value.

Table 3-7. Portals Return Codes: Function Return Codes for the Portals API.

Name	Meaning	Functions
PTL_ABORTED	wait/poll operation was aborted	PtlCTPoll(), PtlCTWait(), PtlEQPoll(), PtlEQWait()
PTL_ARG_INVALID	invalid argument passed	<i>all</i> , except PtlAtomicSync(), PtlFini(), PtlHandlelsEqual(), PtlInit()
PTL_CT_NONE_REACHED	timeout reached before any counting event reached the test	PtlCTPoll()
PTL_EQ_DROPPED	at least one event has been dropped	PtlEQGet(), PtlEQPoll(), PtlEQWait()
PTL_EQ_EMPTY	no events available in an event queue	PtlEQGet(), PtlEQPoll()
PTL_FAIL	error during initialization	PtlInit()
PTL_IGNORED	Logical map set failed	PtlSetMap()
PTL_IN_USE	MD, ME, or LE has pending operations	PtlLEUnlink(), PtlMEUnlink()
PTL_LIST_TOO_LONG	list too long	PtlLEAppend(), PtlMEAppend()
PTL_NO_INIT	uninitialized API	<i>all</i> , except PtlFini(), PtlHandlelsEqual(), PtlInit()
PTL_NO_SPACE	insufficient memory	PtlCTAlloc(), PtlEQAlloc(), PtlGetMap(), PtlLEAppend(), PtlMDBind(), PtlMEAppend(), PtlNlInit(), PtlSetMap(), PtlTriggeredPut(), PtlTriggeredGet(), PtlTriggeredCTInc(), PtlTriggeredCTSet(), PtlTriggeredAtomic(), PtlTriggeredFetchAtomic(), PtlTriggeredSwap()
PTL_OK	success	<i>all</i> , except PtlFini(), PtlHandlelsEqual()
PTL_PID_IN_USE	pid is in use	PtlNlInit()
PTL_PT_EQ_NEEDED	EQ must be attached when flow control is enabled	PtlPTAlloc()
PTL_PT_FULL	portal table is full	PtlPTAlloc()
PTL_PT_IN_USE	portal table index is busy	PtlPTAlloc(), PtlPTFree()

Table 3-8 summarizes the remaining constant values introduced by the Portals API. The first column in this table presents the symbolic name for the constant, the second column gives a brief description of the value, the third column identifies the type for the value, and the fourth column identifies the section in which the constant is introduced or described.

Table 3-8. Portals Constants: Other Constants Defined by the Portals API.

Name	Meaning	Base Type	Definition
PTL_ACK_REQ	request an acknowledgment	ptl_ack_req_t	3.15
PTL_BAND	Compute and return the bitwise AND of the initiator and target value	ptl_op_t	3.15.4

continued on next page

continued from previous page

Name	Meaning	Base Type	Definition
PTL_BOR	Compute and return the bitwise OR of the initiator and target value	ptl_op_t	3.15.4
PTL_BXOR	Compute and return the bitwise XOR of the initiator and target value	ptl_op_t	3.15.4
PTL_COHERENT_ATOMICS	a flag to indicate that the implementation provides atomic operations which are coherent with processor atomic operations	int	3.15.4
PTL_CSWAP	Conditional swap if target and operand equal	ptl_op_t	3.15.4
PTL_CSWAP_GE	Conditional swap if the operand is greater than or equal to the target	ptl_op_t	3.15.4
PTL_CSWAP_GT	Conditional swap if the operand is greater than the target	ptl_op_t	3.15.4
PTL_CSWAP_LE	Conditional swap if the operand is less than or equal to the target	ptl_op_t	3.15.4
PTL_CSWAP_LT	Conditional swap if the operand is less than the target	ptl_op_t	3.15.4
PTL_CSWAP_NE	Conditional swap if the operand and target are not equal	ptl_op_t	3.15.4
PTL_CT_ACK_REQ	request a counting acknowledgment	ptl_ack_req_t	3.15
PTL_CT_NONE	a NULL count handle	ptl_handle_ct_t	3.3.2
PTL_DIFF	Compute the difference between the target initiator	ptl_op_t	3.15.4
PTL_DOUBLE	64-bit floating-point number	ptl_op_t	3.15.4
PTL_DOUBLE_COMPLEX	64-bit floating-point complex number	ptl_op_t	3.15.4
PTL_EQ_NONE	a NULL event queue handle	ptl_handle_eq_t	3.3.2
PTL_EVENT_ACK	acknowledgment event	ptl_event_kind_t	3.13.1
PTL_EVENT_ATOMIC	atomic event	ptl_event_kind_t	3.13.1
PTL_EVENT_ATOMIC_OVERFLOW	atomic overflow event	ptl_event_kind_t	3.13.1
PTL_EVENT_AUTO_FREE	automatic free event	ptl_event_kind_t	3.13.1
PTL_EVENT_AUTO_UNLINK	automatic unlink event	ptl_event_kind_t	3.13.1
PTL_EVENT_FETCH_ATOMIC	fetching atomic event	ptl_event_kind_t	3.13.1
PTL_EVENT_FETCH_ATOMIC_OVERFLOW	fetching atomic overflow event	ptl_event_kind_t	3.13.1
PTL_EVENT_GET	get event	ptl_event_kind_t	3.13.1
PTL_EVENT_GET_OVERFLOW	get overflow event	ptl_event_kind_t	3.13.1
PTL_EVENT_LINK	event generated when a list entry links	ptl_event_kind_t	3.13.1

continued on next page

Name	Meaning	Base Type	Definition
PTL_EVENT_PT_DISABLED	portal table entry disabled event	ptl_event_kind_t	3.13.1
PTL_EVENT_PUT	put event	ptl_event_kind_t	3.13.1
PTL_EVENT_PUT_OVERFLOW	put overflow event	ptl_event_kind_t	3.13.1
PTL_EVENT_REPLY	reply event	ptl_event_kind_t	3.13.1
PTL_EVENT_SEARCH	search event	ptl_event_kind_t	3.13.1
PTL_EVENT_SEND	send event	ptl_event_kind_t	3.13.1
PTL_FLOAT	32-bit floating-point number	ptl_op_t	3.15.4
PTL_FLOAT_COMPLEX	32-bit floating-point complex number	ptl_op_t	3.15.4
PTL_IFACE_DEFAULT	default interface	ptl_interface_t	3.3.5
PTL_INT16_T	16-bit signed integer	ptl_op_t	3.15.4
PTL_INT32_T	32-bit signed integer	ptl_op_t	3.15.4
PTL_INT64_T	64-bit signed integer	ptl_op_t	3.15.4
PTL_INT8_T	8-bit signed integer	ptl_op_t	3.15.4
PTL_INVALID_HANDLE	invalid handle	ptl_handle_any_t	3.3.2
PTL_IOVEC	a flag to enable scatter/gather memory descriptors	int	3.12.1
PTL_LAND	Compute and return the logical AND of the initiator and target	ptl_op_t	3.15.4
PTL_LE_EVENT_COMM_DISABLE	a flag to disable events associated with new communications	int	3.11.1
PTL_LE_EVENT_CT_BYTES	a flag to count bytes instead of operations	int	3.11.1
PTL_LE_EVENT_CT_COMM	a flag to count communication events	int	3.11.1
PTL_LE_EVENT_CT_OVERFLOW	a flag to count overflow events	int	3.11.1
PTL_LE_EVENT_FLOWCTRL_DISABLE	a flag to disable events associated with flow control	int	3.11.1
PTL_LE_EVENT_LINK_DISABLE	a flag to disable link events	int	3.11.1
PTL_LE_EVENT_OVER_DISABLE	a flag to disable overflow events	int	3.11.1
PTL_LE_EVENT_SUCCESS_DISABLE	a flag to disable events that indicate success	int	3.11.1
PTL_LE_EVENT_UNLINK_DISABLE	a flag to disable unlink events	int	3.11.1
PTL_LE_IS_ACCESSIBLE	a flag to indicate the entire LE is accessible	int	3.11.1
PTL_LE_OP_GET	a flag to enable <i>get</i> operations	int	3.11.1
PTL_LE_OP_PUT	a flag to enable <i>put</i> operations	int	3.11.1
PTL_LE_UNEXPECTED_HDR_DISABLE	a flag to disable adding headers to the unexpected headers list	int	3.11.1

continued from previous page

Name	Meaning	Base Type	Definition
PTL_LE_USE_ONCE	a flag to indicate that the list entry will only be used once	int	3.11.1
PTL_LONG_DOUBLE	System defined long double type	ptl_op_t	3.15.4
PTL_LONG_DOUBLE_COMPLEX	System defined long double complex type	ptl_op_t	3.15.4
PTL_LOR	Compute and return the logical OR of the initiator and target	ptl_op_t	3.15.4
PTL_LXOR	Compute and return the logical XOR of the initiator and target	ptl_op_t	3.15.4
PTL_MAX	Compute and return the maximum of the initiator and target	ptl_op_t	3.15.4
PTL_MD_EVENT_CT_ACK	a flag to count acknowledgment events	int	3.10.1
PTL_MD_EVENT_CT_BYTES	a flag to count bytes instead of operations	int	3.10.1
PTL_MD_EVENT_CT_REPLY	a flag to count reply events	int	3.10.1
PTL_MD_EVENT_CT_SEND	a flag to count send events	int	3.10.1
PTL_MD_EVENT_SEND_DISABLE	a flag to disable send events	int	3.10.1
PTL_MD_EVENT_SUCCESS_DISABLE	a flag to disable events that indicate success	int	3.10.1
PTL_MD_UNORDERED	a flag to indicate that messages from this MD do not need to be ordered	int	3.10.1
PTL_MD_UNRELIABLE	a flag to indicate that delivery of messages from this MD need not be guaranteed	int	3.10.1
PTL_MD_VOLATILE	a flag to indicate that the application will modify the put buffer immediately upon operation return, before receiving a send event.	int	3.10.1
PTL_ME_EVENT_COMM_DISABLE	a flag to disable events associated with new communications	int	3.12.1
PTL_ME_EVENT_CT_BYTES	a flag to count bytes instead of operations	int	3.12.1
PTL_ME_EVENT_CT_COMM	a flag to count communication events	int	3.12.1
PTL_ME_EVENT_CT_OVERFLOW	a flag to count overflow events	int	3.12.1
PTL_ME_EVENT_FLOWCTRL_DISABLE	a flag to disable events associated with flow control	int	3.12.1
PTL_ME_EVENT_LINK_DISABLE	a flag to disable link events	int	3.12.1

continued on next page

continued from previous page

Name	Meaning	Base Type	Definition
PTL_ME_EVENT_OVER_DISABLE	a flag to disable overflow events	int	3.12.1
PTL_ME_EVENT_SUCCESS_DISABLE	a flag to disable events that indicate success	int	3.12.1
PTL_ME_EVENT_UNLINK_DISABLE	a flag to disable unlink events	int	3.12.1
PTL_ME_IS_ACCESSIBLE	a flag to indicate the entire ME is accessible	int	3.12.1
PTL_ME_LOCAL_INC_UH_RLENGTH	a flag to indicate that the local offset should be incremented by the requested length of each matching header found in the unexpected headers list	int	3.12.1
PTL_ME_MANAGE_LOCAL	a flag to enable the use of local offsets	int	3.12.1
PTL_ME_MAY_ALIGN	a flag to indicate that the implementation may align an incoming message to a natural boundary to enhance performance	int	3.12.1
PTL_ME_NO_TRUNCATE	a flag to disable truncation of a request	int	3.12.1
PTL_ME_OP_GET	a flag to enable <i>get</i> operations	int	3.12.1
PTL_ME_OP_PUT	a flag to enable <i>put</i> operations	int	3.12.1
PTL_ME_UNEXPECTED_HDR_DISABLE	a flag to disable adding headers to the unexpected headers list	int	3.12.1
PTL_ME_USE_ONCE	a flag to indicate that the match list entry will only be used once	int	3.12.1
PTL_MIN	Compute and return the minimum of the initiator and target	ptl_op_t	3.15.4
PTL_MSWAP	A masked version of the swap operation	ptl_op_t	3.15.4
PTL_NI_DROPPED	message was dropped	ptl_ni_fail_t	3.13.3
PTL_NI_LOGICAL	a flag to indicate that the network interface must provide logical addresses for network endpoints	int	3.6.2
PTL_NI_MATCHING	a flag to indicate that the network interface must provide matching portals addressing	int	3.6.2

continued on next page

continued from previous page

Name	Meaning	Base Type	Definition
PTL_NI_NO_MATCH	search of the unexpected list did not find an entry or all matching entries were found and all corresponding events were generated	ptl_ni_fail_t	3.13
PTL_NI_NO_MATCHING	a flag to indicate that the network interface must provide non-matching portals addressing	int	3.6.2
PTL_NI_OK	successful event	ptl_ni_fail_t	3.13.3
PTL_NI_OP_VIOLATION	message encountered an operation violation	ptl_ni_fail_t	3.13.3
PTL_NI_PERM_VIOLATION	message encountered a permissions violation	ptl_ni_fail_t	3.13.3
PTL_NI_PHYSICAL	a flag to indicate that the network interface must provide physical addresses for network endpoints	int	3.6.2
PTL_NI_PT_DISABLED	message encountered a disabled portal table entry	ptl_ni_fail_t	3.13.3
PTL_NI_SEGV	message attempted to access inaccessible memory	ptl_ni_fail_t	3.13.3
PTL_NI_UNDELIVERABLE	message could not be delivered	ptl_ni_fail_t	3.13.3
PTL_NID_ANY	wildcard for node identifier fields	ptl_nid_t	3.3.6
PTL_NO_ACK_REQ	request no acknowledgment	ptl_ack_req_t	3.15
PTL_OC_ACK_REQ	request an operation completed acknowledgment	ptl_ack_req_t	3.15
PTL_OVERFLOW_LIST	specifies the overflow list attached to a portal table entry	int	3.12.2
PTL_PID_ANY	wildcard for process identifier fields	ptl_pid_t	3.3.6
PTL_PID_MAX	Maximum legal process identifier	ptl_pid_t	3.6.2
PTL_PRIORITY_LIST	specifies the priority list attached to a portal table entry	int	3.12.2
PTL_PROD	Compute and return the product of the initiator and target value	ptl_op_t	3.15.4
PTL_PT_ALLOC_DISABLED	a flag to indicate a portal table entry should be initialized to disabled status	int	3.7.1
PTL_PT_ANY	wildcard for portal table entry identifier fields	ptl_pt_index_t	3.7.1

continued on next page

continued from previous page

Name	Meaning	Base Type	Definition
PTL_PT_FLOWCTRL	a flag to request flow control	int	3.7.1
PTL_PT_ONLY_TRUNCATE	a flag to indicate that the priority list on this portal table entry will only have entries without the PTL_ME_NO_TRUNCATE option set	int	3.7.1
PTL_PT_ONLY_USE_ONCE	a flag to indicate that the priority list on this portal table entry will only have entries with the PTL_ME_USE_ONCE or PTL_LE_USE_ONCE option set	int	3.7.1
PTL_RANK_ANY	wildcard for rank fields	ptl_rank_t	3.3.6
PTL_SEARCH_DELETE	specifies that the unexpected list should be searched and the matching item should be deleted	int	3.12.4
PTL_SEARCH_ONLY	specifies that the unexpected list should only be searched	int	3.12.4
PTL_SIZE_MAX	maximum value of a ptl_size_t	ptl_size_t	3.3.1
PTL_SR_DROP_COUNT	index for the dropped count register	ptl_sr_index_t	3.3.7
PTL_SR_OPERATION_VIOLATIONS	index for the operation violations register	ptl_sr_index_t	3.3.7
PTL_SR_PERMISSION_VIOLATIONS	index for the permission violations register	ptl_sr_index_t	3.3.7
PTL_SUM	Compute and return the sum of the initiator and target	ptl_op_t	3.15.4
PTL_SWAP	Swap the initiator and target value	ptl_op_t	3.15.4
PTL_TARGET_BIND_INACCESSIBLE	A flag to indicate that the implementation should allow LEs to be bound over ranges of memory that are not allocated	int	3.6.1
PTL_TIME_FOREVER	a flag to indicate unbounded time	ptl_time_t	3.13.9
PTL_TOTAL_DATA_ORDERING	A flag to indicate that the implementation should attempt to provide total data ordering	int	3.6.1
PTL_UID_ANY	wildcard for usage identifier	ptl_uid_t	3.3.6
PTL_UINT16_T	16-bit unsigned integer	ptl_op_t	3.15.4
PTL_UINT32_T	32-bit unsigned integer	ptl_op_t	3.15.4

continued on next page

continued from previous page

Name	Meaning	Base Type	Definition
PTL_UINT64_T	64-bit unsigned integer	ptl_op_t	3.15.4
PTL_UINT8_T	8-bit unsigned integer	ptl_op_t	3.15.4

4. Guide to Implementors

In this chapter, we provide a number of notes and clarifications useful to implementors of the Portals specification. This chapter is not normative; that is, this chapter only seeks to clarify and raise subtle points in the standard. Should any statement in this chapter conflict with statements in another chapter, the other chapter is correct.

4.1. Run-time Support

The Portals API does not include a run-time interface; this is assumed to be provided by other sources, such as the machine system software or as part of an upper-layer protocol. This is similar to Open Fabrics, Myrinet/MX, and TCP/IP, which provide communication semantics, but say little about process lifespan or interaction. Interaction with a run-time is clearly unavoidable due to logically addressed network interfaces, but the proper interaction between the run-time and `PtlSetMap()/PtlGetMap()` is the responsibility of the upper layer protocol.

Many implementations of the Portals specification (both Portals 4 and earlier specifications) were tightly coupled with a specific run-time. It is expected that such coupling will continue on tightly integrated platforms in which Portals is the lowest layer communication interface. While the user of the portals library must always call `PtlSetMap()` before using a logically addressed interface, the implementation is free to ignore the requested mapping and provide it's own by returning `PTL_IGNORED`.

4.2. Data Transfer

The Portals API uses five types of messages: *put*, *acknowledgment*, *get*, *reply*, and *atomic*. In this section, we describe the information passed on the wire for each type of message. We also describe how this information is used to process incoming messages. The Portals specification does not enforce a given wire protocol or in what order and what manner information is passed along the communication path.

4.2.1. Sending Messages

Table 4-1 summarizes the information that is transmitted for a *put* request. The first column provides a descriptive name for the information, the second column provides the type for this information, the third column identifies the source of the information, and the fourth column provides additional notes. Most information that is transmitted is obtained directly from the *put* operation.

It may not be necessary for the implementation to transmit all fields listed in Table 4-1. For example, portals semantics require that an *acknowledgment* event contains the *user_ptr* and it must be placed in the event queue referenced by the *eq_handle* found in the MD referenced by the *md_handle* associated with the *put*; i.e., the *acknowledgment* event provides a pointer that the application can use to identify the operation and must be placed in the right memory descriptor's event queue. One approach would be to send the *user_ptr* and *md_handle* to the *target* in the *put* and back again in the *acknowledgment* message. If an implementation has another way of tracking the *user_ptr* and *md_handle* at the initiator, then sending the *user_ptr* and *md_handle* should not be necessary.

Notice that the *match_bits*, *md_handle* and *user_ptr* fields in the *put* operation are optional. If the *put* is originating from a non-matching network interface, there is no need for the *match_bits* to be transmitted since the destination

Table 4-1. Send Request: Information Passed in a Send Request — PtlPut().

Information	Type	PtlPut() Argument	Notes
operation	int		indicates a <i>put</i> request
ack type	ptl_ack_req_t	<i>ack_req</i>	
options	unsigned int	<i>md_handle</i>	<i>options</i> field from NI associated with MD
initiator	ptl_process_t		local information
usage	ptl_uid_t		local information
target	ptl_process_t	<i>target_id</i>	
portal index	ptl_pt_index_t	<i>pt_index</i>	
match bits	ptl_match_bits_t	<i>match_bits</i>	opt. if <i>options</i> .PTL_NI_NO_MATCHING
offset	ptl_size_t	<i>remote_offset</i>	
memory desc	ptl_handle_md_t	<i>md_handle</i>	opt. if <i>ack_req</i> =PTL_NO_ACK_REQ
header data	ptl_hdr_data_t	<i>hdr_data</i>	user data in header
put user pointer	void *	<i>user_ptr</i>	opt. if <i>ack_req</i> =PTL_NO_ACK_REQ or <i>ack_req</i> =PTL_CT_ACK_REQ or <i>ack_req</i> =PTL_OC_ACK_REQ
length	ptl_size_t	<i>length</i>	<i>length</i> argument
data	bytes	<i>md_handle</i>	user data

will ignore them. Similarly, if no acknowledgment was requested, *md_handle* and *user_ptr* do not need to be sent. If an acknowledgment is requested (either PTL_ACK_REQ, PTL_CT_ACK_REQ, or PTL_OC_ACK_REQ), then the *md_handle* may be sent in the *put* message so that the *target* can send it back to the *initiator* in the *acknowledgment* message. The *md_handle* is needed by the *initiator* to find the right event queue for the acknowledgment event. The *user_ptr* is only required in the case of a full acknowledgment (PTL_ACK_REQ). PTL_CT_ACK_REQ and PTL_OC_ACK_REQ requests do not require the *user_ptr* field to generate the acknowledgment event at the *initiator* of the *put* operation.

A portals header contains 8 bytes of user supplied data specified by the *hdr_data* argument passed to PtlPut(). This is useful for out-of-band data transmissions with or without bulk data. The header bytes are stored in the event generated at the *target*. (See Section 3.15.2 on page 94.)

Tables 4-2 and 4-3 summarizes the information transmitted in an *acknowledgment*. Most of the information is simply echoed from the *put* request. Notice that the *initiator* and *target* are obtained directly from the *put* request but are swapped in generating the *acknowledgment*. The only new pieces of information in the *acknowledgment* are the manipulated length, which is determined as the *put* request is satisfied, and the actual offset used.

Table 4-2. Acknowledgment: Information Passed in an Acknowledgment.

Information	Type	PtlPut() Argument	Notes
operation	int		indicates an <i>acknowledgment</i>
options	unsigned int	<i>put_md_handle</i>	<i>options</i> field from NI associated with MD
initiator	ptl_process_t	<i>target_id</i>	echo <i>target</i> of <i>put</i>
target	ptl_process_t	<i>initiator</i>	echo <i>initiator</i> of <i>put</i>
memory descriptor	ptl_handle_md_t	<i>md_handle</i>	echo <i>md_handle</i> of <i>put</i>
put user pointer	void *	<i>user_ptr</i>	echo <i>user_ptr</i> of <i>put</i>
offset	ptl_size_t	<i>remote_offset</i>	obtained from the operation
manipulated length	ptl_size_t		obtained from the operation
matched list	ptl_list_t		obtained from the operation

If an *acknowledgment* has been requested, the associated memory descriptor remains in use by the implementation until the *acknowledgment* arrives and can be logged in the event queue. See Section 3.10.4 for how pending

Table 4-3. Acknowledgment: Information Passed in a “Counting” Acknowledgment.

Information	Type	PtlPut() Argument	Notes
operation	int		indicates an <i>acknowledgment</i>
options	unsigned int	<i>put_md_handle</i>	<i>options</i> field from NI associated with MD
initiator	ptl_process_t	<i>target_id</i>	local information on <i>put</i> target
target	ptl_process_t	<i>initiator</i>	echo <i>initiator</i> of <i>put</i>
memory descriptor	ptl_handle_md_t	<i>md_handle</i>	echo <i>md_handle</i> of <i>put</i>
manipulated length	ptl_size_t		obtained from the operation

operations affect when memory descriptors may be unlinked.

If the target match list entry has the PTL_ME_MANAGE_LOCAL flag set, the offset local to the *target* match list entry is used. If the flag is not set, the offset requested by the *initiator* is used. An *acknowledgment* message returns the actual value used.

Lightweight “counting” acknowledgments do not require the actual offset used or user pointer since they do not generate a **ptl_event_t** at the *put* operation *initiator*.

Table 4-4 summarizes the information that is transmitted for a *get* request. Like the information transmitted in a *put* request, most of the information transmitted in a *get* request is obtained directly from the **PtlGet()** operation. The memory descriptor must not be unlinked until the *reply* is received.

Table 4-4. Get Request: Information Passed in a Get Request — **PtlGet()** and **PtlGetRegion()**.

Information	Type	PtlGet() Argument	Notes
operation	int		indicates a <i>get</i> operation
options	unsigned int	<i>md_handle</i>	<i>options</i> field from NI associated with MD
initiator	ptl_process_t		local information
usage	ptl_uid_t		local information
target	ptl_process_t	<i>target_id</i>	
portal index	ptl_pt_index_t	<i>pt_index</i>	
match bits	ptl_match_bits_t	<i>match_bits</i>	optional if the PTL_NI_NO_MATCHING option is set.
offset	ptl_size_t	<i>remote_offset</i>	
memory descriptor	ptl_handle_md_t	<i>md_handle</i>	destination of <i>reply</i>
length	ptl_size_t	<i>length</i>	
initiator offset	ptl_size_t	<i>local_offset</i>	
get user pointer	void *	<i>user_ptr</i>	

Table 4-5 summarizes the information transmitted in a *reply*. Like an *acknowledgment*, most of the information is simply echoed from the *get* request. The *initiator* and *target* are obtained directly from the *get* request but are swapped in generating the *reply*. The only new information in the *reply* are the manipulated length, the actual offset used, and the data, which are determined as the *get* request is satisfied.

Table 4-6 presents the information that needs to be transmitted from the *initiator* to the *target* for an *atomic* operation. The result of an *atomic* operation is a *reply* and (optionally) an *acknowledgment* as described in Table 4-5.

Table 4-5. Reply: Information Passed in a Reply.

Information	Type	PtlGet() Argument	Notes
operation	int		indicates an <i>reply</i>
options	unsigned int	<i>get_md_handle</i>	<i>options</i> field from NI associated with MD
initiator	ptl_process_t	<i>target_id</i>	local information on <i>get</i> target
target	ptl_process_t	<i>initiator</i>	echo <i>initiator</i> of <i>get</i>
memory descriptor	ptl_handle_md_t	<i>md_handle</i>	echo <i>md_handle</i> of <i>get</i>
initiator offset	ptl_size_t	<i>local_offset</i>	echo <i>local_offset</i> of <i>get</i>
get user pointer	void *	<i>user_ptr</i>	echo <i>user_ptr</i> of <i>get</i>
manipulated length	ptl_size_t		obtained from the operation
offset	ptl_size_t	<i>remote_offset</i>	obtained from the operation
matched list	ptl_list_t		obtained from the operation
data	<i>bytes</i>		obtained from the operation

Table 4-6. Atomic Request: Information Passed in an Atomic Request.

Information	Type	PtlAtomic() Argument	Notes
operation	int		indicates the type of <i>atomic</i> operation and datatype
options	unsigned int	<i>put_md_handle</i>	<i>options</i> field from NI associated with MD
ack type	ptl_ack_req_t	<i>ack_req</i>	
initiator	ptl_process_t		local information
usage	ptl_uid_t		local information
target	ptl_process_t	<i>target_id</i>	
portal index	ptl_pt_index_t	<i>pt_index</i>	
memory descriptor	ptl_handle_md_t	<i>put_md_handle</i>	opt. if <i>ack_req</i> =PTL_NO_ACK_REQ
user pointer	void *	<i>user_ptr</i>	opt. if <i>ack_req</i> =PTL_NO_ACK_REQ or <i>ack_req</i> =PTL_CT_ACK_REQ or <i>ack_req</i> =PTL_OC_ACK_REQ
match bits	ptl_match_bits_t	<i>match_bits</i>	optional if the PTL_NI_NO_MATCHING option is set.
offset	ptl_size_t	<i>remote_offset</i>	
memory descriptor	ptl_handle_md_t	<i>get_md_handle</i>	destination of <i>reply</i>
length	ptl_size_t	<i>put_md_handle</i>	<i>length</i> member
operand	<i>bytes</i>	operand	Used in CSWAP and MSWAP operations
data	<i>bytes</i>	<i>put_md_handle</i>	user data

4.2.2. Receiving Messages

When an incoming message arrives on a network interface, the communication system first checks that the *target* process identified in the request is a valid process that has initialized the network interface (i.e., that the *target* process has a valid portal table). If this test fails, the communication system discards the message and increments the dropped message count for the interface. The remainder of the processing depends on the type of the incoming message. *put*, *get*, and *atomic* messages go through portals address translation (searching a list) and must then pass an access control test. In contrast, *acknowledgment* and *reply* messages bypass the access control checks and the translation step.

Acknowledgment messages include the memory descriptor handle used in the original **PtlPut()** operation. This memory descriptor will identify the event queue where the event should be recorded. Upon receipt of an acknowledgment, the runtime system only needs to confirm that the memory descriptor and event queue still exist. Should any of these conditions fail, the message is simply discarded, and the dropped message count for the interface

is incremented. Otherwise, the system builds an acknowledgment event from the information in the acknowledgment message and adds it to the event queue.

Reception of *reply* messages is also relatively straightforward. Each *reply* message includes a memory descriptor handle. If this descriptor exists, it is used to receive the message. A *reply* message will be dropped if the memory descriptor identified in the request does not exist or it has become inactive. In this case, the dropped message count for the interface is incremented. Every memory descriptor accepts and truncates incoming *reply* messages, eliminating the other potential reasons for rejecting a *reply* message.

The critical step in processing an incoming *put*, *get*, or *atomic* request involves mapping the request to a match list entry (or list entry). This step starts by using the portal index in the incoming request to identify a list of match list entries (or list entries). On a matching interface, the list of match list entries is searched in sequential order until a match list entry is found whose match criteria matches the match bits in the incoming request and that accepts the request. On a non-matching interface, the first item on the list is used and a permissions check is performed.

Because *acknowledgment* and *reply* messages are generated in response to requests made by the process receiving these messages, the checks performed by the runtime system for acknowledgments and replies are minimal. In contrast, *put*, *get*, and *atomic* messages are generated by remote processes and the checks performed for these messages are more extensive. Incoming *put*, *get*, or *atomic* messages may be rejected because:

- the portal index supplied in the request is not valid;
- the match bits supplied in the request do not match any of the match list entries that accepts the request, or
- the access control information provided in the list entry does not match the information provided in the message.

In all cases, if the message is rejected, the incoming message is discarded and the dropped message count for the interface is incremented.

A list entry or match list entry may reject an incoming request if the `PTL_ME_OP_PUT` or `PTL_ME_OP_GET` option has not been enabled and the operation is *put*, *get*, or *atomic* (Table 4-7). In addition, a match list entry may reject an incoming request if the length specified in the request is too long for the match list entry and the `PTL_ME_NO_TRUNCATE` option has been enabled. Truncation is always enabled on standard list entries; thus, a message cannot be rejected for this reason on a non-matching network interface.

Also see Sections 2.5 and Figure 2-11.

Table 4-7. Portals Operations and ME/LE Flags: A - indicates that the operation will be rejected, and a • indicates that the operation will be accepted.

Target ME/LE Flags	Operation				
	<i>put</i>	<i>get</i>	<i>atomic</i>	<code>PtlSwap()</code>	<code>PtlFetchAtomic()</code>
none	-	-	-	-	-
<code>PTL_ME_OP_PUT/PTL_LE_OP_PUT</code>	•	-	•	•	-
<code>PTL_ME_OP_GET/PTL_LE_OP_GET</code>	-	•	-	-	•
both	•	•	•	•	•

4.3. Event Generation and Error Reporting

The types of events and when they are generated is discussed in Chapter 3.13. Operations related to memory descriptors, list entries, and match list entries may both generate a full event (of type `ptl_event_t`) and update a counting event. There is no implied ordering between the generation of a full event and updating of a counting event, although if the user requests both a full event and a counting event, the implementation must deliver both in a timely fashion.

Acknowledgment events require special attention due to the flexibility Portals provides the user in controlling acknowledgments. An acknowledgment event is only generated if the *initiator* requests an acknowledgment and either the *target* enables sending an acknowledgment in the list entry or an error occurs during the operation. Requesting a full acknowledgment (PTL_ACK_REQ) without an event queue on the associated memory descriptor (or with success events disabled) still results in the generation of a counting event.

Bibliography

- [1] N.R. Adiga and et. al. An Overview of the BlueGene/L Supercomputer. In *In Proceedings of the SC 2002 Conference on High Performance Networking and Computing*, Baltimore, MD, November 2002.
- [2] Robert Alverson. Red Storm. In *Invited Talk, Hot Chips 15*, August 2003.
- [3] Christian Bell and Dan Bonachea. A new dma registration strategy for pinning-based high performance networks. *Parallel and Distributed Processing Symposium, International*, 0:198a, 2003.
- [4] Ron Brightwell, David S. Greenberg, Arthur B. Maccabe, and Rolf Riesen. Massively Parallel Computing with Commodity Components. *Parallel Computing*, 26:243–266, February 2000.
- [5] Ron Brightwell, Trammell Hudson, Rolf Riesen, and Arthur B. Maccabe. The Portals 3.0 message passing interface. Technical Report SAND99-2959, Sandia National Laboratories, December 1999.
- [6] Ron Brightwell and Arthur B. Maccabe. Scalability limitations of VIA-based technologies in supporting MPI. In *Fourth MPI Developers' and Users' Conference*, March 2000.
- [7] Ron Brightwell and Lance Shuler. Design and implementation of MPI on Puma portals. In *Proceedings of the Second MPI Developer's Conference*, pages 18–25, July 1996.
- [8] Barbara Chapman, Tony Curtis, Swaroop Pophale, Stephen Poole, Jeff Kuehn, Chuck Koelbel, and Lauren Smith. Introducing openshmem: Shmem for the pgas community. In *Proceedings of the Fourth Conference on Partitioned Global Address Space Programming Model*, page 2. ACM, 2010.
- [9] Compaq, Microsoft, and Intel. Virtual Interface Architecture Specification Version 1.0. Technical report, Compaq, Microsoft, and Intel, December 1997.
- [10] Cray Research, Inc. *SHMEM Technical Note for C, SG-2516 2.3*, October 1994.
- [11] Infiniband Trade Association. <http://www.infinibandta.org>, 1999.
- [12] Y. Ishikawa, H. Tezuka, and A. Hori. PM: A High-Performance Communication Library for Multi-user Parallel Environments. Technical Report TR-96015, RWCP, 1996.
- [13] Mario Lauria, Scott Pakin, and Andrew Chien. Efficient Layering for High Speed Communication: Fast Messages 2.x. In *Proceedings of the IEEE International Symposium on High Performance Distributed Computing*, 1998.
- [14] Arthur B. Maccabe, Kevin S. McCurley, Rolf Riesen, and Stephen R. Wheat. SUNMOS for the Intel Paragon: A brief user's guide. In *Proceedings of the Intel Supercomputer Users' Group. 1994 Annual North America Users' Conference.*, pages 245–251, June 1994.
- [15] Message Passing Interface Forum. MPI: A Message-Passing Interface standard. *The International Journal of Supercomputer Applications and High Performance Computing*, 8:159–416, 1994.
- [16] Message Passing Interface Forum. *MPI: A Message-Passing Interface Version 3.1*, June 2015.
- [17] Myricom, Inc. The GM Message Passing System. Technical report, Myricom, Inc., 1997.
- [18] Rolf Riesen, Ron Brightwell, and Arthur B. Maccabe. The evolution of Portals, an API for high performance communication. *To be published*, 2005.

- [19] Rolf Riesen, Ron Brightwell, Arthur B. Maccabe, Trammell Hudson, and Kevin Pedretti. The Portals 3.3 message passing interface: Document revision 2.0. Technical report SAND2006-0420, Sandia National Laboratories, January 2006.
- [20] Lance Shuler, Chu Jong, Rolf Riesen, David van Dresser, Arthur B. Maccabe, Lee Ann Fisk, and T. Mack Stallcup. The Puma operating system for massively parallel computers. In *Proceeding of the 1995 Intel Supercomputer User's Group Conference*. Intel Supercomputer User's Group, 1995.
- [21] Task Group of Technical Committee T11. Information Technology - Scheduled Transfer Protocol - Working Draft 2.0. Technical report, Accredited Standards Committee NCITS, July 1998.

APPENDIX A. Portals Design Guidelines

Early versions of Portals were based on the idea of using data structures to describe to the transport mechanism how data should be delivered. This worked well for the Puma OS on the Intel Paragon but not so well under Linux on Cplant. The solution was to create a thin API over those data structures and add a level of abstraction. The result was Portals 3.x. While Portals 3.x supported MPI well for kernel level implementations, more advanced offloading network interfaces and the rising importance of PGAS models exposed several weaknesses. This led to several enhancements that became Portals 4.x.

When designing and expanding this API, we were guided by several principles and requirements. We have divided them into three categories: requirements that must be fulfilled by the API and its implementations, requirements that should be met, and a wish list of things that would be nice if Portals 4.x could provide them.

A.1. Mandatory Requirements

Message passing protocols. Portals *must* support efficient implementations of commonly used message passing protocols.

Partitioned Global Address Space (PGAS) Support. Portals *must* support efficient implementations of typical PGAS languages and programming interfaces.

Portability. It *must* be possible to develop implementations of Portals on a variety of existing message passing interfaces.

Scalability. It *must* be possible to write efficient implementations of Portals for systems with millions of nodes.

Performance. It *must* be possible to write high performance (e.g., low latency, high bandwidth) implementations of Portals on existing hardware and on hardware capable of offloading Portals processing.

Multiprocess support. Portals *must* support use of the communication interface by tens of processes per node.

Communication between processes from different executables. Portals *must* support the ability to pass messages between processes instantiated from different executables.

Runtime independence. The ability of a process to perform message passing *must not* depend on the existence of an external runtime environment, scheduling mechanism, or other special utilities outside of normal UNIX process startup.

Memory protection. Portals *must* ensure that a process cannot access the memory of another process without consent.

A.2. The *Will* Requirements

Operational API. Portals *will* be defined by operations, not modifications to data structures. This means that the interface will have explicit operations to send and receive messages. (It does not mean that the receive operation will involve a copy of the message body.)

MPI. It *will* be possible to write an efficient implementation of the point-to-point operations in MPI 1 using Portals.

PGAS. It *will* be possible to write an efficient implementation of the one-sided and atomic operations found in PGAS models using Portals.

Network Interfaces. It *will* be possible to write an efficient implementation of Portals using a network interface that provides offload support.

Operating Systems. It *will* be possible to write an efficient implementation of Portals using a lightweight kernel *or* Linux as the host OS.

Message Size. Portals *will not* impose an arbitrary restriction on the size of message that can be sent.

OS bypass. Portals *will* support an OS bypass message passing strategy. That is, high performance implementations of the message passing mechanisms will be able to bypass the OS and deliver messages directly to the application.

Put/Get. Portals *will* support remote put/get operations.

Packets. It *will* be possible to write efficient implementations of Portals that packetize message transmission.

Receive operation. The receive operation of Portals *will* use an address and length pair to specify where the message body should be placed.

Receiver managed communication. Portals *will* support receive-side management of message space, and this management will be performed during message receipt.

Sender managed communication. Portals *will* support send-side management of message space.

Parallel I/O. Portals *will* be able to serve as the transport mechanism for a parallel file I/O system.

Gateways. It *will* be possible to write *gateway* processes using Portals. A gateway process is a process that receives messages from one implementation of Portals and transmits them to another implementation of Portals.

Asynchronous operations. Portals *will* support asynchronous operations to allow computation and communication to overlap.

Receive side matching. Portals *will* allow matching on the receive side before data is delivered into the user buffer.

A.3. The *Should* Requirements

Message Alignment. Portals *should not* impose any restrictions regarding the alignment of the address(es) used to specify the contents of a message.

Striping. Portals *should* be able to take advantage of multiple interfaces on a single logical network to improve the bandwidth

Socket API. Portals *should* support an efficient implementation of sockets (including UDP and TCP/IP).

Internetwork consistency. Portals *should not* impose any consistency requirements across multiple networks/interfaces. In particular, there will not be any memory consistency/coherency requirements when messages arrive on independent paths.

Ease of use. Programming with Portals *should* be no more complex than programming traditional message passing environments such as UNIX sockets or MPI. An in-depth understanding of the implementation or access to implementation-level information should not be required.

Minimal API. Only the smallest number of functions and definitions necessary to manipulate the data structures should be specified. That means, for example, that convenience functions, which can be implemented with the already defined functions, will not become part of the API.

APPENDIX B. README Definition

Each Portals implementation should provide a README file that details implementation-specific choices. This appendix describes such a file by listing which parameters should be specified.

Limits. The call `PtINIInit()` accepts a desired set of limits and returns a set of actual limits. The README should state the possible ranges of actual limits for this implementation, as well as the acceptable ranges for the values passed into `PtINIInit()`. See Section 3.6.1

Resource Usage. The implementation will be required to consume some user memory for the limits specified in `PtINIInit()`. The README should document the memory resources required by the implementation and should enumerate the relationship between the memory resources consumed and the limits requested in the desired set of limits passed into `PtINIInit()`. See Section 3.6.1

Status Registers. Portals define a set of status registers (Section 3.3.7). The type `ptl_sr_index_t` defines the mandatory `PTL_SR_DROP_COUNT`, `PTL_SR_PERMISSION_VIOLATIONS`, and `PTL_SR_OPERATION_VIOLATIONS`, as well as all other, implementation specific indexes. The README should list what indexes are available and what their purposes are.

Network interfaces. Each Portals implementation defines `PTL_IFACE_DEFAULT` to access the default network interface on a system (Sections 3.3.5 and 3.6.2). An implementation that supports multiple interfaces must specify the constants used to access the various interfaces through `PtINIInit()`.

Portal table. The Portals specification says that a compliant implementation must provide at least 250 entries per portal table (Section 3.6). The README file should state how many entries will actually be provided.

Alignment. If an implementation favors specific alignments for memory descriptors, the README should state what they are and the (performance) consequences if they are not observed (Sections 3.10.1 and 3.12.1). Furthermore, if the implementation supports unaligned atomic operations, it should be documented.

APPENDIX C. Summary of Changes

This chapter documents significant changes to the Portals specification between releases. Semantic changes are always noted and significant corrections to the specification are also noted.

C.1. Portals 4.3

- Added new Section 2.2 and text in Sections 3.10.1 and 3.15.2 providing the new option `PTL_MD_UNRELIABLE` that allows for network dropping of puts when set. **PtlPut()** is only allowed to use the `PTL_MD_UNRELIABLE` option when `PTL_NO_ACK_REQ` is also used.
- Added new function, **PtlAbort()**, and return code, `PTL_ABORTED`. Added new Section 3.5.3 to define the function. Added language to Sections 3.13.6, 3.13.8, 3.13.9, 3.14.3, 3.14.6 and 3.14.7, describing impact on **PtlEQFree()**, **PtlEQWait()**, **PtlEQPoll()**, **PtlCTFree()**, **PtlCTWait()**, **PtlCTPoll()**, respectively.
- Added new option `PTL_PT_ALLOC_DISABLED` to **PtlPTAlloc()** (Section 3.7.1).
- Added new option `PTL_ME_LOCAL_INC_UH_RLENGTH` to the match list entry type, Section 3.12.1. Language in Sections 3.12 and 3.12.2 updated to reflect the new option. For 4.3, `PTL_ME_LOCAL_INC_UH_RLENGTH` is ignored by **PtlMSearch()**, and language making this explicit was added to Section 3.12.4.
- Updated semantics of the `PTL_NI_NO_MATCH` failure constant (used with the `PTL_EVENT_SEARCH` event type) to indicate that either no matching message is found or all matching messages have been returned (and there are no remaining matching messages). Added description of updated semantics to sections 3.11.4 (**PtlESearch()**), 3.12.4 (**PtlMSearch()**), 3.13.1 (description of `PTL_EVENT_SEARCH` event), 3.13.3 (description of `PTL_NI_NO_MATCH` constant), and 3.19 (summary of `PTL_NI_NO_MATCH` constant, Table 3-8).
- Added language to Sections 3.11.4 and 3.12.4 specifying the behavior of counting events when a counting event handle is included in an LE or ME used to conduct a search.
- Deprecated `PTL_INTERRUPTED` as a valid return code (Sections 3.13.8, 3.13.9, 3.14.6, and 3.14.7).
- Deprecated the `PTL_LE_ACK_DISABLE` and `PTL_ME_ACK_DISABLE` options for list entries (Section 3.11.1) and match list entries (Section 3.12.1). Language on failure notification (Section 3.13.3) updated to reflect target-side disabling of acknowledgements is no longer an option, and the options removed from the table of Portals constants (Table 3-8).
- `PTL_NO_SPACE` added as a possible return value for all triggered operations, to indicate the requested triggered operation could not be registered (Sections 3.16.1, 3.16.2, 3.16.3, 3.16.4, 3.16.5, 3.16.6, and 3.16.7).
- The discussion note and list of list entry search operation constants in Section 3.11.4 incorrectly used ‘overflow list’ when ‘unexpected list’ was intended; this has been fixed.
- Added language to Section 3.15.4 clarifying the semantics of atomic operations.

C.2. Portals 4.2

- Updated Section 3.15.4 to ensure that atomics operations that match on the overflow list do not perform atomic operations on the overflow buffer. The atomic operation should be checked for in the event from the overflow list and the host is responsible for ensuring the atomic is applied correctly to the target buffer.
- Clarified the behavior of truncation checks for MEs with the `PTL_ME_NO_TRUNCATE` option set. Zero byte messages will never cause a match to fail, even if their offsets do not pass a truncation check.
- Section 3.12 was updated to clarify that messages with a length of zero bytes are exempt from `PTL_NI_SEGV` error checks.
- Section 3.15.1 was modified to allow for full events to be generated for `PTL_CT_ACK_REQ` and `PTL_OC_ACK_REQ` event types, but does not guarantee that these full events will contain valid entries for all of the fields for a `PTL_EVENT_ACK` shown in Table 3-3.
- Added text in Section 2.5 that emphasizes that overflow events must only be posted after the corresponding payload has been delivered to the appropriate buffer.
- Clarified in Section 2.6 that implementations are allowed to alter data in a user buffer used in a reply operation, but only after that operation has completed.
- Added text in Section 2.7 and 2.7.1 that makes it clear that Portals does not provide write ordering of bytes in individual messages.
- Section 2.7.3 was updated to remove the restriction that “The result of two simultaneous operations targeting the same memory address through different list entries is undefined”. Section 2.7.3 now describes the expected behavior of simultaneous operations targeting the same memory through different list entries (identical to simultaneous operations on the same list entry).
- Text in Section 2.7.4 was re-written to make it clearer that unexpected messages maintain sufficient information to enable MPI matching ordering semantics.
- Section 2.8 removed direct references to InfiniBand operation and clarified that RNR NACKs can inhibit pipelining regardless of network.
- Defined the exact semantics used by Portals for counting event arithmetic in Section 3.14.
- Added a discussion to Section 3.14.9 noting that while it is possible to enter negative values through the C interface for CT calls, they will be interpreted as very large positive numbers.
- Added text in discussion in Section 3.16 encouraging users to order calls to triggered operation calls by increasing trigger threshold.
- Added discussion text for atomic `PTL_DIFF` in Section 3.15.4 to clarify that the operation field in an atomic event may be ignored for successful `PTL_DIFF` operations.
- Updated discussion in Section 3.7.1 to make it clear that the options specified to `PtiPTAlloc()` are binding guarantees consistent with the rest of the text.
- Added text to Section 3.6.1 to clarify that although resource limits are specified separately for each NI type, the resources for each specific NI type are shared amongst all users of that NI.
- Fixed a typo in Section 3.6.1 that erroneously stated the maximum index of PTs was 63, when it should be 249 (250 total PTs).
- Clarified in Section 3.6.2 that in the case where an implementation provided *actual* limits to the NI that are equal to the maximum value of the corresponding type of that limit, the actual resources provided may be bounded by another system resource that is outside of the Portals implementation’s control.

- Added clarification to Section 3.6.6 that an application can only create two logically addressed NIs (one matching and one non-matching). Also noted that the **PTL_IGNORED** return code for **PtlSetMap()** does not indicate success, only that no error occurred and the map was not changed.
- Changed return code for **PtlGetMap()** from **PTL_NO_SPACE** to **PTL_IGNORED** for the case where there was no mapping set for the requested logical network interface as it is not due to lack of memory resources, but a lack of the mapping existing that is the result of the return code (the request was ignored as it was invalid for the requested map).
- Added text to Sections 3.11.4 and 3.12.4 that notes that persistent ME/LEs can cause multiple matches in a **PtlLESearch()** or **PtlMESearch()** call while use-once ME/LEs will only match once.
- Removed discussion text from Sections 3.11.4 and 3.12.4 that was already covered in the preceding **PtlLEAppend()** and **PtlMEAppend()** descriptions in Sections 3.11.2 and 3.12.2.
- Clarified in Section 3.15.2 that **PTL_EVENT_ACK** events indicate that the remote operation was fully processed by the target.
- Added text to Section 3.15.4 to explain what conditions may exist that would result in a locally managed offset that is larger than *mlength*.

C.3. Portals 4.1

- Clarified in Section 3.2 that the **PTL_MINOR_VERSION** constant will be the first integer value in the minor version numbering (e.g. 4.0.2 minor version is 0).
- Changed the name of user ID to usage ID in Section 3.8.1. This is meant to clarify that UIDs do not need to be tied to the system defined user identification values.
- Clarified that for interfaces that support **PTL_COHERENT_ATOMICS** **PtlAtomicSync()** is not required and that atomics are coherent with host operations in 3.15.4.
- Defined protected headers in section 2.5 with regards to usage IDs. Protected headers were referred to in the text in previous Portals versions but never explicitly defined.
- Clarified the role of **PtlAtomicSync()** in 3.15.8 with respect to the visibility of data on the host.
- Provided a guarantee in section 3.16 that operations that trigger on the same threshold value will trigger in the order in which the triggered operations were posted.
- Clarified that in section 2.5 when matching occurs from a **PtlMEAppend()** or **PtlLEAppend()** in the overflow list, the *mlength* returned is that of the original placement due to the overflow list buffer posted, not that copied into the user supplied buffer. Length checking does not occur on overflow list matching and therefore the user must check that the *mlength* is not greater than the length of the posted buffer to be assured that truncation has not occurred.
- Corrected and clarified table 4-7 concerning what operations are accepted based on the target's ME/LE flags.
- Clarified that the options specified to **PtlPTAlloc()** in section 3.7.1 are binding guarantees, not non-binding hints.
- Added **PTL_DIFF** atomic operation in section 3.15.4, which computes the difference between the target and initiator. While **PTL_DIFF** is equivalent to negating the initiator and performing a **PTL_SUM**, **PTL_DIFF** potentially saves a memory copy.
- Increase the minimum number of available portal table entries from 64 to 250 (Section 3.6).
- In section 3.14.6 clarified that the input 'test' in the arguments section is greater than or equal to this value on successful return as stated in the function description.

- Clarify that either local or remote event completions are sufficient conditions to allow the user to reuse a buffer in section 3.13.2.
- Changed `PTL_LE_EVENT_SUCCESS_DISABLE` to only disable success events that are counted by `PTL_LE_EVENT_CT_COMM` and `PTL_LE_EVENT_CT_OVERFLOW` (section 3.11).
- Noted in Section 2.8 that local operations continue to be processed on a disabled portals table entry.

C.4. Portals 4.0.2

- Clarify that `PTL_EVENT_AUTO_FREE` delivery should also consider counting events.
- Note that the event in `PtlEQGet()`, `PtlEQPoll()`, and `PtlEQWait()` is copied from the EQ not removed and that a pointer to the local copy of the `ptl_event_t` structure is returned
- Change the semantics of fairness for `PtlEQPoll()`, events are now returned on the first queue in the array of handles. Rationale: Supporting round robin fairness between different calls of `PtlEQPoll()` and determining when a given array of handles matches a previously used array of handles and tracking where in the round robin rotation each `PtlEQPoll()` call should start at requires significant state tracking by the implementation as well as a comparison of the entire event queue handles array to determine if it has been used previously.
- Change triggered events to only occur on the value of successes rather than successes and failures. This allows for some possible recovery should failures occur and be recoverable.
- Change `PTL_TOTAL_DATA_ORDERING` to provide total data ordering only when both initiator and target NIs have enabled TDO. When targets do not support TDO, they are free to not provide it for incoming messages. Clarify the differences between *max_waw_ordered_size* and `PTL_TOTAL_DATA_ORDERING`.
- Add specification version preprocessor constants to assist application developers in supporting minor revisions along the Portals 4 lifespan. Also require implementations to provide both `portals.h` and `portals4.h` header files.
- Clarify behavior of locally managed offsets and minimum free autounlink for persistent match list entries on the priority list which match headers in the unexpected headers list.
- Clarify that `PTL_LE_EVENT_COMM_DISABLE` and `PTL_ME_EVENT_COMM_DISABLE` should also disable the generation of `PTL_EVENT_FETCH_ATOMIC` events. Also clarify that `PTL_LE_EVENT_CT_COMM` and `PTL_ME_EVENT_CT_COMM` should cause the counting event to be updated for `PTL_EVENT_FETCH_ATOMIC` events.
- Clarified the behavior of `PTL_EVENT_AUTO_FREE` when used with a list entry/match list entry which disables unexpected headers.
- Clarified that `PtlSwap()` returns the value from the target prior to the operation being performed.
- Updated `PtlAtomic()` function descriptions to be consistent with text on `PTL_COHERENT_ATOMICS`.
- Clarify that `PtlEQPoll()` and `PtlCTPoll()` may be called with event queues and counting event handles associated with different logical network interfaces as long as they share the same physical network interface.
- Change the semantics of `PTL_MD_VOLATILE`. Implementations are required to provided volatile semantics if the operation length is less than or equal to *max_volatile_size*, but also must not return an error solely because the operation length is greater than *max_volatile_size*.
- Added `PTL_EVENT_ERROR` event. this event is intended to be used when an unspecified error may be detectable and recoverable by an application.
- Added clarification to atomics *length* argument that if it is not a integral multiple of the datatype size that it may be truncated to zero, a multiple of the datatype or return `PTL_ARG_INVALID`.

- Added PTL_NI_SEGV and PTL_NI_NO_MATCH to the `ptl_ni_fail_t` definition in Section 3.13.3

C.5. Portals 4.0.1

- Specify that PTL_EVENT_AUTO_UNLINK must come after all other events on a list entry / match list entry.
- PTL_EVENT_PT_DISABLED, PTL_EVENT_LINK, PTL_EVENT_AUTO_UNLINK, and PTL_EVENT_AUTO_FREE should provide a `user_ptr` and `ni_fail_type` and not a `start` and `hdr_data` field.
- For clarity, remove “and return” from the description of PTL_MIN, PTL_MAX, PTL_SUM, PTL_PROD, PTL_LOR, PTL_LAND, PTL_BOR, PTL_BAND, PTL_LXOR, PTL_BXOR from Section 3.15.4.

C.6. Portals 4.0

The most recent version of this document described Portals version 3.3 [19]. Since then we have made changes to the API and semantics of Portals, as well as changes to the document. This appendix summarizes the changes between version 3.3 and the current 4.0 version. Many of the fundamental changes were driven by the desire to reduce the tight coupling required between the application processor and the portals processor, but some additions were made to better support lighter weight communications models such as PGAS.

Foremost, Portals version 4.0 was substantially enhanced to better support the various PGAS programming models. Communication operations that do not include matching were added along with key atomic operations. In addition, the ordering definition was substantially strengthened relative to Portals version 3.3 for small messages. In support of the lightweight communication semantics required by PGAS models, lightweight “counting” events and acknowledgments were added. A `PtlAtomic()` function was added to support functionality commonly provided in PGAS models. Finally, the Portals ordering model was substantially expanded to better support some PGAS models.

An equally fundamental change in Portals version 4.0 adds a mechanism to cope better with the concept of unexpected messages in MPI. Whereas version 3.3 used `PtlMDUpdate()` to atomically insert items into the match list so that the MPI implementation could manage unexpected messages, version 4.0 adds an overflow list where the application provides buffer space that the implementation can use to store unexpected messages. The implementation is then responsible for matching new list insertions to items that have arrived and are resident in the overflow list space. This change was necessary to eliminate round trips between the processor and the NIC for each item that was added to the match list (now named the priority list).

A third major change separated all resources for initiators and targets. Memory descriptors are used by the initiator to describe memory regions while list entries are used by targets to describe the memory region *and* matching criteria (in the case of match list entries). This separation of resources was also extended to events, where the number of event types was significantly reduced and only required fields for a given event type must be defined.

To better offload collective operations, a set of *triggered* operations were added. These operations allow an application to build non-blocking, offloaded collective operations with independent progress. They include variants of both the data movement operations (get and put) as well as the atomic operations.

Another set of changes arise from a desire to simplify hardware implementations. The threshold value was removed from the target and was replaced by the ability to specify that a match list entry is “use once” or “persistent”. List insertions occur *only* at the tail of the list, since unexpected message handling has been separated out into a separate list.

Access control entries were found to be a non-scalable resource, so they have been eliminated. At the same time, it was recognized that the PTL_LE_OP_PUT and PTL_LE_OP_GET semantics required a form of matching. These two options along with the ability to include usage ID based authentication were moved to *permissions fields* on the respective list entry or match list entry.

Ordering only at the message level was found to be insufficient for many PGAS models, which often require ordering of data. Unfortunately, uniformly requiring data ordering could create unnecessary performance constraints. As such, the ordering definition has been expanded to include data ordering and to let the user disable that ordering and message ordering.

Index

- ack_req (field), 94, 100, 105, 107, 127, 129
- acknowledgment, *see* operations
- acknowledgment type, **92**, 92
- actual (field), 32, 42–44, 58, 65, 138
- actual_map_size (field), 46, 47
- address translation, 21
- address space opening, 21
- address translation, 23, **26**, 30, 129
- addressing, portals, 34
- alignment, 136
- API, 14, [15]
- API summary, **114**
- application bypass, 18
- application space, 22
- application bypass, **19**, 20, 21
- argument names, *see* structure fields
- ASC, [15]
- ASCII, [15]
- Atomic
 - alignment, 96, 136
- atomic, *see* operations
 - datatypes, 98
 - operations, 97
- atomic operation, 21, 23, 95, 116, 117
- atomic swap, *see* swap
- atomic_operation (field), 81
- atomic_type (field), 81
- background, **18**
- buffer alignment, 136
- buffer alignment, 54, 59, 67
- bypass
 - application, 18, **19**, 20, 21
 - OS, 18, **19**, 20, 135
- CAF, 17
- changes, API and document, 137
- communication model, **19**
- connection-oriented, 17
- connectionless, 17, 19
- constants, 35
 - PTL_ACK_REQ, 35, 92, 93, 118, 127, 131
 - PTL_BAND, 97, 99, 118, 141
 - PTL_BOR, 97, 99, 119, 141
 - PTL_BXOR, 97, 99, 119, 141
 - PTL_COHERENT_ATOMICS, 42, 96, 103, 119, 139, 140
 - PTL_CSWAP, 96, 97, 99, 102, 103, 119
 - PTL_CSWAP_GE, 98, 99, 119
 - PTL_CSWAP_GT, 98, 99, 119
 - PTL_CSWAP_LE, 97, 99, 119
 - PTL_CSWAP_LT, 98, 99, 119
 - PTL_CSWAP_NE, 97, 99, 119
 - PTL_CT_ACK_REQ, 92, 93, 119, 127, 129, 138
 - PTL_CT_NONE, 36, 55, 59, 67, 119
 - PTL_DIFF, 97, 99, 119, 138, 139
 - PTL_DOUBLE, 98, 119
 - PTL_DOUBLE_COMPLEX, 98, 119
 - PTL_EQ_NONE, 36, 48, 55, 77, 119
 - PTL_EVENT_ACK, 34, 54, 57–59, 66, 67, 75, 77, 78, 80, 81, 92–94, 96, 99, 119, 138, 139
 - PTL_EVENT_ATOMIC, 60, 69, 74, 75, 77, 80, 81, 96, 98, 119
 - PTL_EVENT_ATOMIC_OVERFLOW, 60, 61, 63, 69, 72, 74, 77, 81, 96, 98, 119
 - PTL_EVENT_AUTO_FREE, 58, 60, 65, 66, 68, 69, 75, 77, 81, 119, 140, 141
 - PTL_EVENT_AUTO_UNLINK, 58, 60, 65, 66, 68, 69, 75, 77, 81, 119, 141
 - PTL_EVENT_ERROR, 75, 77, 81, 140
 - PTL_EVENT_FETCH_ATOMIC, 60, 69, 74, 75, 77, 80, 81, 100, 119, 140
 - PTL_EVENT_FETCH_ATOMIC_OVERFLOW, 60, 61, 63, 69, 72, 74, 77, 81, 100, 119
 - PTL_EVENT_GET, 60, 69, 74, 75, 77, 80, 81, 94, 119
 - PTL_EVENT_GET_OVERFLOW, 60, 61, 63, 69, 72, 74, 77, 81, 94, 119
 - PTL_EVENT_LINK, 60, 61, 69, 70, 75, 77, 81, 119, 141
 - PTL_EVENT_PT_DISABLED, 28, 34, 49, 60, 61, 69, 70, 75, 77, 81, 82, 120, 141
 - PTL_EVENT_PUT, 27, 60, 61, 69, 70, 74, 75, 77, 80, 81, 93, 98, 120
 - PTL_EVENT_PUT_OVERFLOW, 27, 60, 61, 63, 69, 72, 74, 77, 81, 93, 120
 - PTL_EVENT_REPLY, 34, 54, 57, 75, 77, 78, 81, 94, 96, 100, 102, 120
 - PTL_EVENT_SEARCH, 60, 63, 64, 69, 72, 75, 77, 80, 81, 120, 137
 - PTL_EVENT_SEND, 54, 57, 75, 77, 78, 80, 81,

93, 94, 96, 99, 100, 120
 PTL_FLOAT, 98, 120
 PTL_FLOAT_COMPLEX, 98, 120
 PTL_IFACE_DEFAULT, 37, 120, 136
 PTL_INT16_T, 98, 120
 PTL_INT32_T, 98, 120
 PTL_INT64_T, 98, 120
 PTL_INT8_T, 98, 120
 PTL_INVALID_HANDLE, 36, 113, 120
 PTL_IOVEC, 53, 55, 58, 60, 65, 69, 96, 120
 PTL_LAND, 97, 99, 120, 141
 PTL_LE_ACK_DISABLE, 137
 PTL_LE_EVENT_COMM_DISABLE, 60, 120, 140
 PTL_LE_EVENT_CT_BYTES, 61, 120
 PTL_LE_EVENT_CT_COMM, 60, 120, 140
 PTL_LE_EVENT_CT_OVERFLOW, 61, 120, 140
 PTL_LE_EVENT_FLOWCTRL_DISABLE, 60, 61, 120
 PTL_LE_EVENT_LINK_DISABLE, 60, 120
 PTL_LE_EVENT_OVER_DISABLE, 60, 120
 PTL_LE_EVENT_SUCCESS_DISABLE, 60, 120, 140
 PTL_LE_EVENT_UNLINK_DISABLE, 60, 61, 120
 PTL_LE_IS_ACCESSIBLE, 58, 60, 120
 PTL_LE_OP_GET, 59, 120, 130, 141
 PTL_LE_OP_PUT, 59, 120, 130, 141
 PTL_LE_UNEXPECTED_HDR_DISABLE, 60, 120
 PTL_LE_USE_ONCE, 48, 58, 60, 61, 63, 64, 75, 121, 124
 PTL_LONG_DOUBLE, 96, 98, 121
 PTL_LONG_DOUBLE_COMPLEX, 96, 98, 121
 PTL_LOR, 97, 99, 121, 141
 PTL_LXOR, 97, 99, 121, 141
 PTL_MAJOR_VERSION, 36
 PTL_MAX, 97, 99, 121, 141
 PTL_MD_EVENT_CT_ACK, 54, 121
 PTL_MD_EVENT_CT_BYTES, 55, 92, 121
 PTL_MD_EVENT_CT_REPLY, 54, 121
 PTL_MD_EVENT_CT_SEND, 54, 121
 PTL_MD_EVENT_SEND_DISABLE, 54, 121
 PTL_MD_EVENT_SUCCESS_DISABLE, 54, 121
 PTL_MD_UNORDERED, 31, 33, 55, 121
 PTL_MD_UNRELIABLE, 24, 55, 94, 121, 137
 PTL_MD_VOLATILE, 42, 55, 121, 140
 PTL_ME_ACK_DISABLE, 137
 PTL_ME_EVENT_COMM_DISABLE, 69, 77, 121, 140
 PTL_ME_EVENT_CT_BYTES, 69, 121
 PTL_ME_EVENT_CT_COMM, 69, 121, 140
 PTL_ME_EVENT_CT_OVERFLOW, 69, 121
 PTL_ME_EVENT_FLOWCTRL_DISABLE, 69, 70, 121
 PTL_ME_EVENT_LINK_DISABLE, 69, 121
 PTL_ME_EVENT_OVER_DISABLE, 69, 122
 PTL_ME_EVENT_SUCCESS_DISABLE, 69, 122
 PTL_ME_EVENT_UNLINK_DISABLE, 69, 70, 77, 122
 PTL_ME_IS_ACCESSIBLE, 65, 69, 122
 PTL_ME_LOCAL_INC_UH_RLENGTH, 66, 68, 70, 72, 122, 137
 PTL_ME_MANAGE_LOCAL, 67, 68, 94, 95, 100, 101, 103, 122, 128
 PTL_ME_MAY_ALIGN, 68, 70, 96, 122
 PTL_ME_NO_TRUNCATE, 29, 48, 68, 122, 124, 130, 138
 PTL_ME_OP_GET, 67, 96, 122, 130
 PTL_ME_OP_PUT, 67, 96, 122, 130
 PTL_ME_UNEXPECTED_HDR_DISABLE, 68, 122
 PTL_ME_USE_ONCE, 48, 68, 70, 72, 73, 75, 122, 124
 PTL_MIN, 97, 99, 122, 141
 PTL_MINOR_VERSION, 36, 139
 PTL_MSWAP, 96, 98, 99, 102, 103, 122
 PTL_NI_DROPPED, 78, 122
 PTL_NI_LOGICAL, 37, 42, 43, 122
 PTL_NI_MATCHING, 42, 43, 122
 PTL_NI_NO_MATCH, 64, 72, 75, 78, 123, 137, 141
 PTL_NI_NO_MATCHING, 42, 43, 58, 123, 127-129
 PTL_NI_OK, 63, 72, 75, 77, 81, 123
 PTL_NI_OP_VIOLATION, 59, 67, 78, 123
 PTL_NI_PERM_VIOLATION, 59, 67, 78, 123
 PTL_NI_PHYSICAL, 37, 42, 43, 123
 PTL_NI_PT_DISABLED, 34, 78, 123
 PTL_NI_SEGV, 53, 58, 65, 78, 123, 138, 141
 PTL_NI_UNDELIVERABLE, 78, 123
 PTL_NID_ANY, 37, 69, 123
 PTL_NO_ACK_REQ, 93, 94, 123, 127, 129, 137
 PTL_OC_ACK_REQ, 92, 93, 123, 127, 129, 138
 PTL_OVERFLOW_LIST, 61, 70, 123
 PTL_PID_ANY, 37, 43, 69, 123
 PTL_PID_MAX, 43, 123

PTL_PRIORITY_LIST, 61, 70, 123
 PTL_PROD, 97, 99, 123, 141
 PTL_PT_ALLOC_DISABLED, 47–49, 123, 137
 PTL_PT_ANY, 48, 123
 PTL_PT_FLOWCTRL, 33, 34, 48, 124
 PTL_PT_ONLY_TRUNCATE, 48, 49, 124
 PTL_PT_ONLY_USE_ONCE, 48, 49, 124
 PTL_RANK_ANY, 37, 69, 124
 PTL_SEARCH_DELETE, 63, 64, 72, 124
 PTL_SEARCH_ONLY, 63, 64, 72, 124
 PTL_SIZE_MAX, 36, 53, 58, 65, 124
 PTL_SR_DROP_COUNT, 28, 33, 38, 124, 136
 PTL_SR_OPERATION_VIOLATIONS, 38, 59, 67, 124, 136
 PTL_SR_PERMISSION_VIOLATIONS, 38, 59, 67, 124, 136
 PTL_SUM, 97, 99, 124, 139, 141
 PTL_SWAP, 97, 99, 102, 103, 124
 PTL_TARGET_BIND_INACCESSIBLE, 42, 58, 65, 124
 PTL_TIME_FOREVER, 85, 90, 124
 PTL_TOTAL_DATA_ORDERING, 32, 33, 42, 124, 140
 PTL_UID_ANY, 37, 59, 67, 124
 PTL_UINT16_T, 98, 124
 PTL_UINT32_T, 98, 124
 PTL_UINT64_T, 98, 125
 PTL_UINT8_T, 98, 125
 summary, 118
 count (field), 81, 82
 counting event
 type, **86**
 counting event
 allocate, 86
 enable, 54, 55, 60, 61, 69
 freeing, 87
 freeing triggered operations, 88
 get, 88
 increment, 91
 poll, 90
 set, 91
 triggered increment, 110
 triggered set, 111
 wait, 89
 counting events, **85**, 85
 Cplant, 14
 CPU interrupts, 19
 ct_handle (field), 55, 57, 59, 67, 86–93, 110, 111
 ct_handles (field), 90
 Data Buffers, 30
 data movement, **21**, 26, 34, **92**
 data types, 36, 114
 datatype (field), 99–103, 107, 108, 110
 Deferred Communication Operations, **111**
 end bundle, 113
 limit bundling, 112
 start bundle, 112
 design guidelines, 134
 desired (field), 32, 43
 discarded messages, 129
 discarded events, 93
 discarded messages, 19, 21, 130
 DMA, [15]
 dropped message count, 124, 129, 130
 dropped messages, 38, 83–85, 118
 eq_handle (field), 48, 55, 57, 82–84, 126
 eq_handles (field), 84, 85
 event, 20, 26, 59, **74**
 disable, 60, 69, 120–122
 occurrence, 76
 overflow list, 60, 69
 types, 74, 76, 77
 types (diagram), 76
 unlink, 60, 69
 event (field), 83–85, 89, 90
 event queue, [15]
 allocation, 81
 freeing, 82
 get, 83
 poll, 84
 type, **78**
 wait, 84
 failure (field), 86, 89, 91
 failure notification, 77
 faults, **20**
 features (field), 32, 42, 58, 65, 96
 fetch and atomic operation, 117
 flow control, 81, 82
 support, 33
 user-level, 18
 function return codes, *see* return codes
 functions
 PtlAbort, 38, **39**, 39, 82, 84, 85, 87, 89, 90, 116, 137
 PtlAtomic, 23, 92, 95, 96, 98, **99**, 99, 100, 104, 106, 107, 114–116, 129, 140, 141
 PtlAtomicSync, 96, 103, **104**, 116, 118, 139
 PtlCTAlloc, 34, 86, **87**, 114–116, 118
 PtlCTCancelTriggered, **88**, 88, 104, 114, 117
 PtlCTFree, 34, 39, 86, **87**, 87, 88, 114, 117, 137
 PtlCTGet, 39, 86, **88**, 88, 91, 114, 117
 PtlCTInc, 86, **91**, 91, 104, 110, 114, 117

PtlCTPoll, 34, 39, 86, 87, **90**, 90, 114, 116–118, 137, 140
 PtlCTSet, 86, 88, **91**, 91, 111, 114, 117
 PtlCTWait, 34, 39, 86, 87, **89**, 89, 90, 114, 116–118, 137
 PtlEndBundle, 112, **113**, 113, 115, 117
 PtlEQAlloc, 34, 35, 74, 81, **82**, 115–118
 PtlEQFree, 34, 39, 74, **82**, 82, 115, 117, 137
 PtlEQGet, 39, 74, **83**, 83, 84, 114, 115, 117, 118, 140
 PtlEQPoll, 34, 39, 74, 82–84, **85**, 114–118, 137, 140
 PtlEQWait, 34, 39, 74, 82, 83, **84**, 84, 114, 115, 117, 118, 137, 140
 PtlFetchAtomic, 23, 60, 67, 74, 92, 95, 96, 99, 100, **101**, 102, 108, 114–117, 130
 PtlFini, 38, **39**, 39, 117, 118
 PtlGet, 92, 94, **95**, 96, 106, 115–117, 128, 129
 PtlGetId, 42, **52**, 52, 115–117
 PtlGetMap, 45, 46, **47**, 115–118, 126, 139
 PtlGetPhysId, 42, 45, 52, **53**, 115–117
 PtlGetUid, **51**, 51, 115–117
 PtlHandleIsEqual, **113**, 113, 114, 117, 118
 PtlInit, 35, 38, **39**, 39, 117, 118
 PtlLEAppend, 27, 33, 34, 57, **61**, 61–63, 70, 74–76, 80, 115–118, 139
 PtlLESearch, 63, **64**, 64, 73, 75, 115–117, 137, 139
 PtlLEUnlink, 27, 57, 62, **63**, 63, 115, 117, 118
 PtlMDBind, 37, 53, **56**, 56, 57, 115, 117, 118
 PtlMDRelease, 37, 53, **57**, 57, 115, 117
 PtlMEAppend, 27, 33, 65, **70**, 70–72, 74–76, 80, 115–118, 139
 PtlMESearch, 72, **73**, 73, 75, 115–117, 137, 139
 PtlMEUnlink, 27, 34, 65, 71, **72**, 72, 115, 117, 118
 PtlNIFini, 40, 42, **44**, 44, 115, 117
 PtlNIHandle, 36, 40, **45**, 45, 114, 115, 117
 PtlNIInit, 32, 33, 40, **42**, 42–45, 96, 115, 117, 118, 136
 PtlNIStatus, 37, 40, **44**, 44, 115–117
 PtlPTAlloc, 47, **48**, 82, 115–118, 137–139
 PtlPTDisable, 34, 49, **50**, 78, 115–117
 PtlPTEnable, 34, 47, 49, **50**, 50, 51, 115–117
 PtlPTFree, **49**, 49, 115–118
 PtlPut, 92, **93**, 93, 95, 98, 100, 101, 103, 105, 114–117, 127–129, 137
 PtlSetMap, 42, 45, **46**, 46, 115–118, 126, 139
 PtlStartBundle, **112**, 112, 113, 115, 117
 PtlSwap, 60, 67, 74, 92, 96, 98, 99, **102**, 102, 109, 110, 114–117, 130, 140
 PtlTriggeredAtomic, 104, 106, **107**, 114–118

PtlTriggeredCTInc, 106, **110**, 110, 114, 116–118
 PtlTriggeredCTSet, 107, **111**, 114, 116–118
 PtlTriggeredFetchAtomic, **108**, 108, 114–118
 PtlTriggeredGet, 104, **106**, 106, 114–118
 PtlTriggeredPut, **105**, 105, 110, 114–118
 PtlTriggeredSwap, **109**, 109, 114–118
 summary, 116

gather/scatter, *see* scatter/gather

get, *see* operations

get ID, 52, 53

Get Map, 46

get uid, 51

get_md_handle (field), 96, 100–102, 108, 109, 129

handle, 36

 comparison, 113

 operations, **113**

handle (field), 45

handle1 (field), 113

handle2 (field), 113

hdr_data (field), 76, 80, 94, 100, 101, 103, 105, 107, 108, 110, 127, 141

header data, 94, 115, 127

header, trusted, 51

I/O vector, *see* scatter/gather, 55

ID, **37**

 get, 52, 53

 network interface, 37

 node, *see* node ID

 process, *see* process ID

 thread, *see* thread ID

 uid (get), 51

 usage, *see* usage ID

id (field), 52, 53

identifier, *see* ID

iface (field), 42, 43

ignore bits, 30, 69

ignore_bits (field), 69

implementation notes, 13

implementation, quality, 43

increment (field), 91, 92, 110

indexes, portal, 37

initialization, **38**

initiator, *see also* target, [15], 19, 21, 23, 26, 29, 30, 56, 58, 66, 74–78, 80, 92–96, 100, 127–129, 131

initiator (field), 80

interrupt, 19

interrupt latency, 19

iov_base (field), 56

iov_len (field), 56

LE, **57**

- access control, 26
- alignment, 59
- append, 61
- list types, 61
- options, 59
- pending operation, 62
- permissions, 26
- persistent, 61
- protection, 26
- search, 63, 64
- search and delete, 64
- search operations, 64
- unlink, 57, 62, 75, 117

le (field), 62, 64

le_handle (field), 62, 63

length (field), 42, 53–55, 58, 59, 65, 67, 94–96, 100, 101, 103, 105–109, 127–129, 140

lightweight events, **85**

Limit Usage of Bundling, 112

limits, **40**, 40, 115, 136

Linux, 135

list, [15], **57**

list entries, 19

list entry, *see* LE, 57, 58

local offset, *see* offset

local_get_offset (field), 100–102, 108, 109

local_offset (field), 94, 95, 99, 105–107, 128, 129

local_put_offset (field), 100–103, 108, 109

map_size (field), 46, 47

mapping (field), 46, 47

match bits, 29, 30, 35, 37, 69, 94, 95, 100, 101, 103, 115, 127–130

match ID checking, 71

match list, **65**

match list entry, *see* ME, 58, 65, 69

match_bits (field), 69, 80, 94, 95, 100, 101, 103, 105–109, 126–129

match_id (field), 66, 69, 71

matching address translation, 31

max_atomic_size (field), 30, 41, 96, 99–101, 103

max_cts (field), 41

max_entries (field), 41

max_eqs (field), 41

max_fetch_atomic_size (field), 41, 96, 100, 102

max_iovecs (field), 41

max_list_size (field), 41

max_mds (field), 41

max_msg_size (field), 41

max_pt_index (field), 41

max_triggered_ops (field), 41

max_unexpected_headers (field), 34, 41, 44

max_volatile_size (field), 42, 55, 140

max_war_ordered_size (field), 32, 33, 42

max_waw_ordered_size (field), 32, 33, 41, 140

MD, **53**

- alignment, 54, 136
- bind, 56
- options, 54
- pending operation, 57, 128
- release, 53, 57, 117
- unlink, 128
- unreliable, 55
- volatile, 42, 55

md (field), 56

md_handle (field), 56, 57, 93–95, 99, 105–107, 126–129

ME, **65**

- access control, 30
- access control, 26
- alignment, 67, 136
- append, 70
- free, 77, 119
- ignore bits, *see* ignore bits
- link, 77
- match bits, *see* match bits
- message reject, 130
- options, 67
- pending operation, 72
- permissions, 26, 30
- persistent, 70
- protection, 26, 30
- search, 64, 72, 75, 77, 120
- search and delete, 64
- search operations, 64
- truncate, 68, 122, 130
- unlink, 30, 65, 67, 71, 72, 74, 75, 77, 117, 119

me (field), 71, 73

me_handle (field), 71, 72

memory descriptor, *see also* MD, [15], **53**

message, [15]

message rejection, 130

message operation, [15]

messages, receiving, **129**

messages, sending, **126**

min_free (field), 30, 67, 68, 70

mlength (field), 27, 55, 61, 69, 76, 80, 92, 96, 139

MPI, [15], 17, 18, 26, 62, 71, 94, 134

- progress rule, 18, 20

MPI scalability, 18

MPP, [15]

NAL, [15]

naming conventions, **35**

network, [15]

network independence, 18

network interface, *see also* NI, 19, 35–37, **40**, 42, 58, 129

network interface initialization, **42**

network interfaces
multiple, 136

network scalability, 17

new_ct (field), 91, 111

NI
options, 42
retrieving logical maps, 46
setting logical maps, 45

NI fini, 44

NI handle, 45

NI init, 42

NI status, 44

ni_fail_type (field), 34, 53, 54, 58–60, 63, 65, 67, 69, 72, 75, 77, 81, 141

ni_handle (field), 42–53, 56, 61–64, 70–73, 82, 87, 96, 112, 113

nid (field), 52

node, [15]

node ID, 37

node ID, 26, 30, 37, 51

non-matching address translation, 29

NULL LE, 58

NULL ME, 65

offset, 26, 127–129
local, 67, 68, 80, 93, 95, 122
remote, 67, 80, 94, 95, 100, 101, 103

one-sided operation, 19

opening into address space, 21

operand, 129

operand (field), 96, 102, 103, 110

operation (field), 100, 101, 103, 107, 108, 110

operation violations count, 124

operations
acknowledgment, 23, 44, 74–76, 126–130
atomic, 15, 21, 23, 30, 32, 41, 42, 44, 55, 74–76, 92, **98**, 98, 126, 128–130
atomic sync, 103
atomics, **95**
fetch and atomic, **100**
get, 15, 21, 23, 24, 30, 32–34, 41, 42, 44, 55, 59, 60, 67, 68, 74–77, **94**, 95, 96, 117, 120, 122, 126, 128–130
one-sided, 19
put, 15, 19, 21, 23–25, 30, 32–34, 41, 42, 44, 53, 55, 58–60, 62, 65, 67, 68, 72, 74–77, 92, **93**, 93, 94, 96, 117, 120, 122, 126–130
reply, 23, 30, 41, 44, 75, 77, 95, 126, 128–130
swap, **102**
two-sided, 19, 29
options (field), 42, 48, 54, 59, 67, 127–129

Ordering, 30
adaptive, 33
long messages, 32
overlapping regions, 32
short messages, 32
unexpected messages, 33

ordering semantics, 19, 32, 55

OS bypass, 18, **19**, 20, 135

overflow list, 22, 27, 30, 34, 57, 58, 63, 65, 72, 76, 80, 141

parallel job, 19

pending operation, *see* MD

performance, 134

permission violations count, 124

PGAS, 17, 135

pid (field), 42, 43, 52

portability, 40

portal
indexes, 37
table, 40, 136
table index, 48–50, 57, 65, 127–130

portal table entry, 35, **47**
allocation, 47
disable, 49
enable, 50
freeing, 49

portal table entry disabled event, 120

Portals
early versions, 14
Version 2.0, 14
Version 3.0, 14

portals
addressing, *see* address translation
constants, *see* constants, 35
constants summary, 118
data buffers, **30**
data types, **36**, 114
design, 134
functions, *see* functions
functions summary, 116
handle, 36
multi-threading, 34
naming conventions, 35
operations, *see* operations
ordering, **30**
return codes, *see* return codes
return codes summary, 118
scalability, 19
sizes, 36

portals.h, 35

portals4.h, 35

priority list, [15], 22, 27, 28, 57, 58, 63, 72

process, [15], 34

process ID, 26, 29, 30, 37, 43, **51**, 51–53, 65, 69, 71,
 94, 100, 101, 103, 115
 well known, 43
 progress, 20
 progress rule, 18, 20
 protected space, 22, 23
 PT
 options, 48
 pt_index (field), 48–50, 61–64, 70–73, 80, 94, 95,
 100, 101, 103, 105–109, 127–129
 pt_index_req (field), 48
 PtlAbort (func), 38, **39**, 39, 82, 84, 85, 87, 89, 90,
 116, 137
 PtlAtomic (func), 23, 92, 95, 96, 98, **99**, 99, 100, 104,
 106, 107, 114–116, 129, 140, 141
 PtlAtomicSync (func), 96, 103, **104**, 116, 118, 139
 PtlCTAlloc (func), 34, 86, **87**, 114–116, 118
 PtlCTCancelTriggered (func), **88**, 88, 104, 114, 117
 PtlCTFree (func), 34, 39, 86, **87**, 87, 88, 114, 117,
 137
 PtlCTGet (func), 39, 86, **88**, 88, 91, 114, 117
 PtlCTInc (func), 86, **91**, 91, 104, 110, 114, 117
 PtlCTPoll (func), 34, 39, 86, 87, **90**, 90, 114,
 116–118, 137, 140
 PtlCTSet (func), 86, 88, **91**, 91, 111, 114, 117
 PtlCTWait (func), 34, 39, 86, 87, **89**, 89, 90, 114,
 116–118, 137
 PtlEndBundle (func), 112, **113**, 113, 115, 117
 PtlEQAlloc (func), 34, 35, 74, 81, **82**, 115–118
 PtlEQFree (func), 34, 39, 74, **82**, 82, 115, 117, 137
 PtlEQGet (func), 39, 74, **83**, 83, 84, 114, 115, 117,
 118, 140
 PtlEQPoll (func), 34, 39, 74, 82–84, **85**, 114–118,
 137, 140
 PtlEQWait (func), 34, 39, 74, 82, 83, **84**, 84, 114, 115,
 117, 118, 137, 140
 PtlFetchAtomic (func), 23, 60, 67, 74, 92, 95, 96, 99,
 100, **101**, 102, 108, 114–117, 130
 PtlFini (func), 38, **39**, 39, 117, 118
 PtlGet (func), 92, 94, **95**, 96, 106, 115–117, 128,
 129
 PtlGetId (func), 42, **52**, 52, 115–117
 PtlGetMap (func), 45, 46, **47**, 115–118, 126, 139
 PtlGetPhysId (func), 42, 45, 52, **53**, 115–117
 PtlGetUid (func), **51**, 51, 115–117
 PtlHandleIsEqual (func), **113**, 113, 114, 117, 118
 PtlInit (func), 35, 38, **39**, 39, 117, 118
 PtlLEAppend (func), 27, 33, 34, 57, **61**, 61–63, 70,
 74–76, 80, 115–118, 139
 PtlLESearch (func), 63, **64**, 64, 73, 75, 115–117, 137,
 139
 PtlLEUnlink (func), 27, 57, 62, **63**, 63, 115, 117,
 118
 PtlMDBind (func), 37, 53, **56**, 56, 57, 115, 117, 118
 PtlMDRelease (func), 37, 53, **57**, 57, 115, 117
 PtlMEAppend (func), 27, 33, 65, **70**, 70–72, 74–76,
 80, 115–118, 139
 PtlMESearch (func), 72, **73**, 73, 75, 115–117, 137,
 139
 PtlMEUnlink (func), 27, 34, 65, 71, **72**, 72, 115, 117,
 118
 PtlNIFini (func), 40, 42, **44**, 44, 115, 117
 PtlNIHandle (func), 36, 40, **45**, 45, 114, 115, 117
 PtlNIInit (func), 32, 33, 40, **42**, 42–45, 96, 115, 117,
 118, 136
 PtlNIStatus (func), 37, 40, **44**, 44, 115–117
 PtlPTAlloc (func), 47, **48**, 82, 115–118, 137–139
 PtlPTDisable (func), 34, 49, **50**, 78, 115–117
 PtlPTEnable (func), 34, 47, 49, **50**, 50, 51, 115–117
 PtlPTFree (func), **49**, 49, 115–118
 PtlPut (func), 92, **93**, 93, 95, 98, 100, 101, 103, 105,
 114–117, 127–129, 137
 PtlSetMap (func), 42, 45, **46**, 46, 115–118, 126,
 139
 PtlStartBundle (func), **112**, 112, 113, 115, 117
 PtlSwap (func), 60, 67, 74, 92, 96, 98, 99, **102**, 102,
 109, 110, 114–117, 130, 140
 PtlTriggeredAtomic (func), 104, 106, **107**, 114–118
 PtlTriggeredCTInc (func), 106, **110**, 110, 114,
 116–118
 PtlTriggeredCTSet (func), 107, **111**, 114, 116–118
 PtlTriggeredFetchAtomic (func), **108**, 108,
 114–118
 PtlTriggeredGet (func), 104, **106**, 106, 114–118
 PtlTriggeredPut (func), **105**, 105, 110, 114–118
 PtlTriggeredSwap (func), **109**, 109, 114–118
 PTL_ABORTED (return code), 39, 84, 85, 89, 90,
 118, 137
 PTL_ACK_REQ (const), 35, 92, 93, 118, 127, 131
 PTL_ARG_INVALID (return code), 38, 43–53, 56,
 57, 62, 63, 65, 71–73, 82–85, 87–92, 94–96,
 100, 102, 103, 105, 106, 108–113, 118,
 140
 PTL_BAND (const), 97, 99, 118, 141
 PTL_BOR (const), 97, 99, 119, 141
 PTL_BXOR (const), 97, 99, 119, 141
 PTL_COHERENT_ATOMICS (const), 42, 96, 103,
 119, 139, 140
 PTL_CSWAP (const), 96, 97, 99, 102, 103, 119
 PTL_CSWAP_GE (const), 98, 99, 119
 PTL_CSWAP_GT (const), 98, 99, 119
 PTL_CSWAP_LE (const), 97, 99, 119
 PTL_CSWAP_LT (const), 98, 99, 119
 PTL_CSWAP_NE (const), 97, 99, 119
 PTL_CT_ACK_REQ (const), 92, 93, 119, 127, 129,
 138
 PTL_CT_NONE (const), 36, 55, 59, 67, 119
 PTL_CT_NONE_REACHED (return code), 90,

118

PTL_DIFF (const), 97, 99, 119, 138, 139

PTL_DOUBLE (const), 98, 119

PTL_DOUBLE_COMPLEX (const), 98, 119

PTL_EQ_DROPPED (return code), 83–85, 118

PTL_EQ_EMPTY (return code), 83, 85, 118

PTL_EQ_NONE (const), 36, 48, 55, 77, 119

PTL_EVENT_ACK (const), 34, 54, 57–59, 66, 67, 75, 77, 78, 80, 81, 92–94, 96, 99, 119, 138, 139

PTL_EVENT_ATOMIC (const), 60, 69, 74, 75, 77, 80, 81, 96, 98, 119

PTL_EVENT_ATOMIC_OVERFLOW (const), 60, 61, 63, 69, 72, 74, 77, 81, 96, 98, 119

PTL_EVENT_AUTO_FREE (const), 58, 60, 65, 66, 68, 69, 75, 77, 81, 119, 140, 141

PTL_EVENT_AUTO_UNLINK (const), 58, 60, 65, 66, 68, 69, 75, 77, 81, 119, 141

PTL_EVENT_ERROR (const), 75, 77, 81, 140

PTL_EVENT_FETCH_ATOMIC (const), 60, 69, 74, 75, 77, 80, 81, 100, 119, 140

PTL_EVENT_FETCH_ATOMIC_OVERFLOW (const), 60, 61, 63, 69, 72, 74, 77, 81, 100, 119

PTL_EVENT_GET (const), 60, 69, 74, 75, 77, 80, 81, 94, 119

PTL_EVENT_GET_OVERFLOW (const), 60, 61, 63, 69, 72, 74, 77, 81, 94, 119

PTL_EVENT_LINK (const), 60, 61, 69, 70, 75, 77, 81, 119, 141

PTL_EVENT_PT_DISABLED (const), 28, 34, 49, 60, 61, 69, 70, 75, 77, 81, 82, 120, 141

PTL_EVENT_PUT (const), 27, 60, 61, 69, 70, 74, 75, 77, 80, 81, 93, 98, 120

PTL_EVENT_PUT_OVERFLOW (const), 27, 60, 61, 63, 69, 72, 74, 77, 81, 93, 120

PTL_EVENT_REPLY (const), 34, 54, 57, 75, 77, 78, 81, 94, 96, 100, 102, 120

PTL_EVENT_SEARCH (const), 60, 63, 64, 69, 72, 75, 77, 80, 81, 120, 137

PTL_EVENT_SEND (const), 54, 57, 75, 77, 78, 80, 81, 93, 94, 96, 99, 100, 120

PTL_FAIL (return code), 39, 118

PTL_FLOAT (const), 98, 120

PTL_FLOAT_COMPLEX (const), 98, 120

PTL_IFACE_DEFAULT (const), 37, 120, 136

PTL_IGNORED (return code), 46, 47, 118, 126, 139

PTL_IN_USE (return code), 62, 63, 71, 72, 118

PTL_INT16_T (const), 98, 120

PTL_INT32_T (const), 98, 120

PTL_INT64_T (const), 98, 120

PTL_INT8_T (const), 98, 120

PTL_INVALID_HANDLE (const), 36, 113, 120

PTL_IOVEC (const), 53, 55, 58, 60, 65, 69, 96, 120

PTL_LAND (const), 97, 99, 120, 141

PTL_LE_ACK_DISABLE (const), 137

PTL_LE_EVENT_COMM_DISABLE (const), 60, 120, 140

PTL_LE_EVENT_CT_BYTES (const), 61, 120

PTL_LE_EVENT_CT_COMM (const), 60, 120, 140

PTL_LE_EVENT_CT_OVERFLOW (const), 61, 120, 140

PTL_LE_EVENT_FLOWCTRL_DISABLE (const), 60, 61, 120

PTL_LE_EVENT_LINK_DISABLE (const), 60, 120

PTL_LE_EVENT_OVER_DISABLE (const), 60, 120

PTL_LE_EVENT_SUCCESS_DISABLE (const), 60, 120, 140

PTL_LE_EVENT_UNLINK_DISABLE (const), 60, 61, 120

PTL_LE_IS_ACCESSIBLE (const), 58, 60, 120

PTL_LE_OP_GET (const), 59, 120, 130, 141

PTL_LE_OP_PUT (const), 59, 120, 130, 141

PTL_LE_UNEXPECTED_HDR_DISABLE (const), 60, 120

PTL_LE_USE_ONCE (const), 48, 58, 60, 61, 63, 64, 75, 121, 124

PTL_LIST_TOO_LONG (return code), 62, 71, 118

PTL_LONG_DOUBLE (const), 96, 98, 121

PTL_LONG_DOUBLE_COMPLEX (const), 96, 98, 121

PTL_LOR (const), 97, 99, 121, 141

PTL_LXOR (const), 97, 99, 121, 141

PTL_MAJOR_VERSION (const), 36

PTL_MAX (const), 97, 99, 121, 141

PTL_MD_EVENT_CT_ACK (const), 54, 121

PTL_MD_EVENT_CT_BYTES (const), 55, 92, 121

PTL_MD_EVENT_CT_REPLY (const), 54, 121

PTL_MD_EVENT_CT_SEND (const), 54, 121

PTL_MD_EVENT_SEND_DISABLE (const), 54, 121

PTL_MD_EVENT_SUCCESS_DISABLE (const), 54, 121

PTL_MD_UNORDERED (const), 31, 33, 55, 121

PTL_MD_UNRELIABLE (const), 24, 55, 94, 121, 137

PTL_MD_VOLATILE (const), 42, 55, 121, 140

PTL_ME_ACK_DISABLE (const), 137

PTL_ME_EVENT_COMM_DISABLE (const), 69, 77, 121, 140

PTL_ME_EVENT_CT_BYTES (const), 69, 121

PTL_ME_EVENT_CT_COMM (const), 69, 121, 140

PTL_ME_EVENT_CT_OVERFLOW (const), 69, 121
 PTL_ME_EVENT_FLOWCTRL_DISABLE (const), 69, 70, 121
 PTL_ME_EVENT_LINK_DISABLE (const), 69, 121
 PTL_ME_EVENT_OVER_DISABLE (const), 69, 122
 PTL_ME_EVENT_SUCCESS_DISABLE (const), 69, 122
 PTL_ME_EVENT_UNLINK_DISABLE (const), 69, 70, 77, 122
 PTL_ME_IS_ACCESSIBLE (const), 65, 69, 122
 PTL_ME_LOCAL_INC_UH_RLENGTH (const), 66, 68, 70, 72, 122, 137
 PTL_ME_MANAGE_LOCAL (const), 67, 68, 94, 95, 100, 101, 103, 122, 128
 PTL_ME_MAY_ALIGN (const), 68, 70, 96, 122
 PTL_ME_NO_TRUNCATE (const), 29, 48, 68, 122, 124, 130, 138
 PTL_ME_OP_GET (const), 67, 96, 122, 130
 PTL_ME_OP_PUT (const), 67, 96, 122, 130
 PTL_ME_UNEXPECTED_HDR_DISABLE (const), 68, 122
 PTL_ME_USE_ONCE (const), 48, 68, 70, 72, 73, 75, 122, 124
 PTL_MIN (const), 97, 99, 122, 141
 PTL_MINOR_VERSION (const), 36, 139
 PTL_MSWAP (const), 96, 98, 99, 102, 103, 122
 PTL_NI_DROPPED (const), 78, 122
 PTL_NI_LOGICAL (const), 37, 42, 43, 122
 PTL_NI_MATCHING (const), 42, 43, 122
 PTL_NI_NO_MATCH (const), 64, 72, 75, 78, 123, 137, 141
 PTL_NI_NO_MATCHING (const), 42, 43, 58, 123, 127–129
 PTL_NI_OK (const), 63, 72, 75, 77, 81, 123
 PTL_NI_OP_VIOLATION (const), 59, 67, 78, 123
 PTL_NI_PERM_VIOLATION (const), 59, 67, 78, 123
 PTL_NI_PHYSICAL (const), 37, 42, 43, 123
 PTL_NI_PT_DISABLED (const), 34, 78, 123
 PTL_NI_SEGV (const), 53, 58, 65, 78, 123, 138, 141
 PTL_NI_UNDELIVERABLE (const), 78, 123
 PTL_NID_ANY (const), 37, 69, 123
 PTL_NO_ACK_REQ (const), 93, 94, 123, 127, 129, 137
 PTL_NO_INIT (return code), 43–53, 56, 57, 62, 63, 65, 71–73, 82–85, 87–92, 94, 95, 100, 102–107, 109–113, 118
 PTL_NO_SPACE (return code), 43, 46, 56, 62, 71, 82, 87, 105–107, 109–111, 118, 137, 139
 PTL_OC_ACK_REQ (const), 92, 93, 123, 127, 129, 138
 PTL_OK (return code), 35, 39, 43–54, 56, 57, 60, 62, 63, 65, 69, 71–73, 82–85, 87–92, 94, 95, 100, 102–107, 109–113, 118
 PTL_OVERFLOW_LIST (const), 61, 70, 123
 PTL_PID_ANY (const), 37, 43, 69, 123
 PTL_PID_IN_USE (return code), 43, 118
 PTL_PID_MAX (const), 43, 123
 PTL_PRIORITY_LIST (const), 61, 70, 123
 PTL_PROD (const), 97, 99, 123, 141
 PTL_PT_ALLOC_DISABLED (const), 47–49, 123, 137
 PTL_PT_ANY (const), 48, 123
 PTL_PT_EQ_NEEDED (return code), 48, 118
 PTL_PT_FLOWCTRL (const), 33, 34, 48, 124
 PTL_PT_FULL (return code), 48, 118
 PTL_PT_IN_USE (return code), 48, 49, 118
 PTL_PT_ONLY_TRUNCATE (const), 48, 49, 124
 PTL_PT_ONLY_USE_ONCE (const), 48, 49, 124
 PTL_RANK_ANY (const), 37, 69, 124
 PTL_SEARCH_DELETE (const), 63, 64, 72, 124
 PTL_SEARCH_ONLY (const), 63, 64, 72, 124
 PTL_SIZE_MAX (const), 36, 53, 58, 65, 124
 PTL_SR_DROP_COUNT (const), 28, 33, 38, 124, 136
 PTL_SR_OPERATION_VIOLATIONS (const), 38, 59, 67, 124, 136
 PTL_SR_PERMISSION_VIOLATIONS (const), 38, 59, 67, 124, 136
 PTL_SUM (const), 97, 99, 124, 139, 141
 PTL_SWAP (const), 97, 99, 102, 103, 124
 PTL_TARGET_BIND_INACCESSIBLE (const), 42, 58, 65, 124
 PTL_TIME_FOREVER (const), 85, 90, 124
 PTL_TOTAL_DATA_ORDERING (const), 32, 33, 42, 124, 140
 PTL_UID_ANY (const), 37, 59, 67, 124
 PTL_UINT16_T (const), 98, 124
 PTL_UINT32_T (const), 98, 124
 PTL_UINT64_T (const), 98, 125
 PTL_UINT8_T (const), 98, 125
 ptl_ack_req_t (type), 92, 114, 118, 119, 123, 127, 129
 ptl_ct_event_t (type), 65, 86, 89, 90, 114, 116
 ptl_datatype_t (type), 96, 98–100, 102, 114
 ptl_event_kind_t (type), 74, 114, 119, 120
 ptl_event_t (type), 74, 78, 79, 81, 83–85, 114–116, 128, 130, 140
 ptl_handle_any_t (type), 36, 114, 120
 ptl_handle_ct_t (type), 36, 85, 86, 114, 119
 ptl_handle_eq_t (type), 36, 74, 115, 119
 ptl_handle_le_t (type), 115
 ptl_handle_md_t (type), 115, 127–129
 ptl_handle_me_t (type), 115

- ptl_handle_ni_t (type), 36, 115
- ptl_hdr_data_t (type), 115, 127
- ptl_interface_t (type), 37, 115, 120
- PTL_INTERRUPTED (return code), 137
- ptl_iovec_t (type), 53, 55, 56, 58, 60, 65, 69, 115, 116
- ptl_le_t (type), 59, 114–116
- ptl_list (field), 61, 62, 70, 71, 80
- ptl_list_t (type), 61, 70, 115, 127, 129
- ptl_match_bits_t (type), 35, 37, 115, 127–129
- ptl_md_t (type), 31, 33, 54, 114–116
- ptl_me_t (type), 66, 114–116
- ptl_ni_fail_t (type), 77, 115, 122, 123, 141
- ptl_ni_limits_t (type), 33, 40, 115, 116
- ptl_nid_t (type), 37, 115, 123
- ptl_op_t (type), 96–98, 115, 118–125
- ptl_pid_t (type), 37, 115, 123
- ptl_process_t (type), 46, 47, 52, 69, 115, 116, 127–129
- ptl_pt_index_t (type), 37, 116, 123, 127–129
- ptl_rank_t (type), 37, 116, 124
- ptl_search_op (field), 63, 65, 72, 73
- ptl_search_op_t (type), 64, 116
- ptl_size_t (type), 36, 116, 124, 127–129
- ptl_sr_index_t (type), 37, 38, 116, 124, 136
- ptl_sr_value_t (type), 38, 116
- ptl_time_t (type), 116, 124
- ptl_uid_t (type), 37, 116, 124, 127–129
- Puma, 18
- purpose, **17**
- put, *see* operations
- put_md_handle (field), 96, 100–103, 108, 109, 127–129

- quality implementation, 43
- quality of implementation, 19

- rank, 22, 30, 37, 45, 46, 51–53
- rank (field), 37, 46, 52
- README, 35, 136
- receiver-managed, 18
- reliable communication, 22
- remote offset, *see* offset
- remote_offset (field), 80, 94, 95, 100, 101, 103, 105–109, 127–129
- reply, *see* operations
- return codes, **38**, 118
 - PTL_ABORTED, 39, 84, 85, 89, 90, 118, 137
 - PTL_ARG_INVALID, 38, 43–53, 56, 57, 62, 63, 65, 71–73, 82–85, 87–92, 94–96, 100, 102, 103, 105, 106, 108–113, 118, 140
 - PTL_CT_NONE_REACHED, 90, 118
 - PTL_EQ_DROPPED, 83–85, 118
 - PTL_EQ_EMPTY, 83, 85, 118
 - PTL_FAIL, 39, 118
 - PTL_IGNORED, 46, 47, 118, 126, 139
 - PTL_IN_USE, 62, 63, 71, 72, 118
 - PTL_LIST_TOO_LONG, 62, 71, 118
 - PTL_NO_INIT, 43–53, 56, 57, 62, 63, 65, 71–73, 82–85, 87–92, 94, 95, 100, 102–107, 109–113, 118
 - PTL_NO_SPACE, 43, 46, 56, 62, 71, 82, 87, 105–107, 109–111, 118, 137, 139
 - PTL_OK, 35, 39, 43–54, 56, 57, 60, 62, 63, 65, 69, 71–73, 82–85, 87–92, 94, 95, 100, 102–107, 109–113, 118
 - PTL_PID_IN_USE, 43, 118
 - PTL_PT_EQ_NEEDED, 48, 118
 - PTL_PT_FULL, 48, 118
 - PTL_PT_IN_USE, 48, 49, 118
 - PTL_INTERRUPTED, 137
 - summary, 118
- rlength (field), 27, 76, 80
- RMPP, [16]

- scalability, **19**, 134
 - guarantee, 19
 - MPI, 18
 - network, 17
- scatter/gather, 54–56, 59, 60, 67, 69, 115, 120
- Search
 - event generation, 64, 73
 - status registers, 64, 73
- send, 21
- send event, 93, 96, 100, 120
- Set Map, 45
- SHMEM, 17
 - shmem_fence(), 32
- size (field), 85, 90
- sizes, 36
- space
 - application, 22
 - protected, 22
- split event sequence, *see* event start/end
- start (field), 53–55, 58, 59, 65, 67, 76, 80, 141
- state, 19
- status registers, 136
- status (field), 45
- status registers, 37
- status_register (field), 45
- structure fields and argument names
 - ack_req, 94, 100, 105, 107, 127, 129
 - actual, 32, 42–44, 58, 65, 138
 - actual_map_size, 46, 47
 - atomic_operation, 81
 - atomic_type, 81
 - count, 81, 82
 - ct_handle, 55, 57, 59, 67, 86–93, 110, 111

ct_handles, 90
 datatype, 99–103, 107, 108, 110
 desired, 32, 43
 eq_handle, 48, 55, 57, 82–84, 126
 eq_handles, 84, 85
 event, 83–85, 89, 90
 failure, 86, 89, 91
 features, 32, 42, 58, 65, 96
 get_md_handle, 96, 100–102, 108, 109, 129
 handle, 45
 handle1, 113
 handle2, 113
 hdr_data, 76, 80, 94, 100, 101, 103, 105, 107, 108, 110, 127, 141
 id, 52, 53
 iface, 42, 43
 ignore_bits, 69
 increment, 91, 92, 110
 initiator, 80
 iov_base, 56
 iov_len, 56
 le, 62, 64
 le_handle, 62, 63
 length, 42, 53–55, 58, 59, 65, 67, 94–96, 100, 101, 103, 105–109, 127–129, 140
 local_get_offset, 100–102, 108, 109
 local_offset, 94, 95, 99, 105–107, 128, 129
 local_put_offset, 100–103, 108, 109
 map_size, 46, 47
 mapping, 46, 47
 match_bits, 69, 80, 94, 95, 100, 101, 103, 105–109, 126–129
 match_id, 66, 69, 71
 max_atomic_size, 30, 41, 96, 99–101, 103
 max_cts, 41
 max_entries, 41
 max_eqs, 41
 max_fetch_atomic_size, 41, 96, 100, 102
 max_iovecs, 41
 max_list_size, 41
 max_mds, 41
 max_msg_size, 41
 max_pt_index, 41
 max_triggered_ops, 41
 max_unexpected_headers, 34, 41, 44
 max_volatile_size, 42, 55, 140
 max_war_ordered_size, 32, 33, 42
 max_waw_ordered_size, 32, 33, 41, 140
 md, 56
 md_handle, 56, 57, 93–95, 99, 105–107, 126–129
 me, 71, 73
 me_handle, 71, 72
 min_free, 30, 67, 68, 70
 mlength, 27, 55, 61, 69, 76, 80, 92, 96, 139
 new_ct, 91, 111
 ni_fail_type, 34, 53, 54, 58–60, 63, 65, 67, 69, 72, 75, 77, 81, 141
 ni_handle, 42–53, 56, 61–64, 70–73, 82, 87, 96, 112, 113
 nid, 52
 operand, 96, 102, 103, 110
 operation, 100, 101, 103, 107, 108, 110
 options, 42, 48, 54, 59, 67, 127–129
 pid, 42, 43, 52
 pt_index, 48–50, 61–64, 70–73, 80, 94, 95, 100, 101, 103, 105–109, 127–129
 pt_index_req, 48
 ptl_list, 61, 62, 70, 71, 80
 ptl_search_op, 63, 65, 72, 73
 put_md_handle, 96, 100–103, 108, 109, 127–129
 rank, 37, 46, 52
 remote_offset, 80, 94, 95, 100, 101, 103, 105–109, 127–129
 rlength, 27, 76, 80
 size, 85, 90
 start, 53–55, 58, 59, 65, 67, 76, 80, 141
 status, 45
 status_register, 45
 success, 86, 89, 91
 target_id, 94, 95, 100, 101, 103, 105–109, 127–129
 test, 89, 90
 tests, 90
 threshold, 104–108, 110, 111
 timeout, 84, 85, 90
 trig_ct_handle, 104–108, 110, 111
 type, 80
 uid, 51, 59, 67, 80
 user_ptr, 62, 65, 71, 73, 75, 80, 81, 94, 95, 100, 101, 103, 105–109, 126–129, 141
 which, 84, 85, 90
 success (field), 86, 89, 91
 summary, **114**
 SUNMOS, [16], 18
 swap operation, 117
 target, *see also* initiator, 15, [16], 19, 21, 23, 26, 51, 74–78, 92–96, 100, 101, 103, 126–129, 131
 target_id (field), 94, 95, 100, 101, 103, 105–109, 127–129
 TCP/IP, 17, 135
 test (field), 89, 90
 tests (field), 90
 thread, [16], 34
 thread ID, 52

threshold (field), 104–108, 110, 111
 timeout, 84, 90
 timeout (field), 84, 85, 90
 trig_ct_handle (field), 104–108, 110, 111
 triggered operations, 32, **104**
 atomic, **106**
 canceling, 88
 counting event increment, 110
 counting event set, 111
 fetch and atomic, **108**
 get, **106**
 put, **105**
 swap, **109**
 threshold, 104
 truncate, 68, 122, 130
 trusted header, 51
 two-sided operation, 19, 29
 type (field), 80
 types, *see* data types
 ptl_ack_req_t, 92, 114, 118, 119, 123, 127, 129
 ptl_ct_event_t, 65, 86, 89, 90, 114, 116
 ptl_datatype_t, 96, 98–100, 102, 114
 ptl_event_kind_t, 74, 114, 119, 120
 ptl_event_t, 74, 78, 79, 81, 83–85, 114–116, 128, 130, 140
 ptl_handle_any_t, 36, 114, 120
 ptl_handle_ct_t, 36, 85, 86, 114, 119
 ptl_handle_eq_t, 36, 74, 115, 119
 ptl_handle_le_t, 115
 ptl_handle_md_t, 115, 127–129
 ptl_handle_me_t, 115
 ptl_handle_ni_t, 36, 115
 ptl_hdr_data_t, 115, 127
 ptl_interface_t, 37, 115, 120
 ptl_iovec_t, 53, 55, 56, 58, 60, 65, 69, 115, 116
 ptl_le_t, 59, 114–116
 ptl_list_t, 61, 70, 115, 127, 129
 ptl_match_bits_t, 35, 37, 115, 127–129
 ptl_md_t, 31, 33, 54, 114–116
 ptl_me_t, 66, 114–116
 ptl_ni_fail_t, 77, 115, 122, 123, 141
 ptl_ni_limits_t, 33, 40, 115, 116
 ptl_nid_t, 37, 115, 123
 ptl_op_t, 96–98, 115, 118–125
 ptl_pid_t, 37, 115, 123
 ptl_process_t, 46, 47, 52, 69, 115, 116, 127–129
 ptl_pt_index_t, 37, 116, 123, 127–129
 ptl_rank_t, 37, 116, 124
 ptl_search_op_t, 64, 116
 ptl_size_t, 36, 116, 124, 127–129
 ptl_sr_index_t, 37, 38, 116, 124, 136
 ptl_sr_value_t, 38, 116
 ptl_time_t, 116, 124
 ptl_uid_t, 37, 116, 124, 127–129
 uid (field), 51, 59, 67, 80
 undefined behavior, 38, 39, 44
 unexpected list, 27, 58, 61, 63, 70, 72
 unexpected message event, 74
 unexpected messages, 18
 unlink, 67
 ME, *see* ME
 unreliable datagram, 24
 UPC, 17
 usage, **24**
 usage ID, 37, **51**, 80, 116, 117, 124
 user data, 62, 65, 71, 73, 94
 user memory, 30
 user space, 19
 user-level bypass, *see* application bypass
 user_ptr (field), 62, 65, 71, 73, 75, 80, 81, 94, 95, 100, 101, 103, 105–109, 126–129, 141
 VIA, [16]
 which (field), 84, 85, 90
 wire protocol, 21, 22, 126
 zero copy, **19**
 zero-length buffer, 58, 65

DISTRIBUTION

Email—Internal (encrypt for OUO)

Name	Org.	Sandia Email Address
Technical Library	1911	sanddocs@sandia.gov



**Sandia
National
Laboratories**

Sandia National Laboratories is a multimission laboratory managed and operated by National Technology & Engineering Solutions of Sandia LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.